

Single studies using the SelfControlledCaseSeries package

Martijn J. Schuemie, Marc A. Suchard and Patrick Ryan

2015-11-25

Contents

1	Introduction	1
2	Installation instructions	2
3	Overview	2
4	Studies with a single drug	2
4.1	Configuring the connection to the server	2
4.2	Preparing the health outcome of interest	3
4.3	Extracting the data from the server	4
4.4	Using a simple model	5
4.5	Adding a pre-exposure window	6
4.6	Splitting risk windows	7
4.7	Including age and seasonality	8
4.8	Considering event-dependent observation time	11
5	Studies with more than one drug	12
5.1	Adding a class of drugs	12
5.2	Adding all drugs	14
6	Acknowledgments	16

1 Introduction

This vignette describes how you can use the `SelfControlledCaseSeries` package to perform a single Self-Controlled Case Series (SCCS) study. We will walk through all the steps needed to perform an exemplar study, and we have selected the well-studied topic of the effect of NSAIDs on gastrointestinal (GI) bleeding-related hospitalization. For simplicity, we focus on one NSAID: diclofenac.

2 Installation instructions

Before installing the `SelfControlledCaseSeries` package make sure you have Java available. Java can be downloaded from www.java.com. For Windows users, RTools is also necessary. RTools can be downloaded from CRAN.

The `SelfControlledCaseSeries` package is currently maintained in a [Github repository](#), and has dependencies on other packages in Github. All of these packages can be downloaded and installed from within R using the `devtools` package:

```
install.packages("devtools")
library(devtools)
install_github("ohdsi/OhdsiRTools")
install_github("ohdsi/SqlRender")
install_github("ohdsi/DatabaseConnector")
install_github("ohdsi/Cyclops")
install_github("ohdsi/SelfControlledCaseSeries")
```

Once installed, you can type `library(SelfControlledCaseSeries)` to load the package.

3 Overview

In the `SelfControlledCaseSeries` package a study requires at least three steps:

1. Loading the necessary data from the database.
2. Transforming the data into a format suitable for an SCCS study. This step includes the creation of covariates based on the variables extracted from the database, such as defining risk windows based on exposures.
3. Fitting the model using conditional Poisson regression.

In the following sections these steps will be demonstrated for increasingly complex studies.

4 Studies with a single drug

4.1 Configuring the connection to the server

We need to tell R how to connect to the server where the data are. `SelfControlledCaseSeries` uses the `DatabaseConnector` package, which provides the `createConnectionDetails` function. Type `?createConnectionDetails` for the specific settings required for the various database management systems (DBMS). For example, one might connect to a PostgreSQL database using this code:

```
connectionDetails <- createConnectionDetails(dbms = "postgresql",
                                             server = "localhost/ohdsi",
                                             user = "joe",
                                             password = "supersecret")

cdmDatabaseSchema <- "my_cdm_data"
cohortDatabaseSchema <- "my_results"
cdmVersion <- "4"
```

The last three lines define the `cdmDatabaseSchema` and `cohortDatabaseSchema` variables, as well as the CDM version. We'll use these later to tell R where the data in CDM format live, where we have stored our cohorts of interest, and what version CDM is used. Note that for Microsoft SQL Server, databaseschemas need to specify both the database and the schema, so for example `cdmDatabaseSchema <- "my_cdm_data.dbo"`.

4.2 Preparing the health outcome of interest

We need to define the outcome for our study. One way to do this is by writing SQL statements against the OMOP CDM that populate a table of events in which we are interested. The resulting table should have the same structure as the `cohort` table in the CDM. For CDM v4, this means it should have the fields `cohort_concept_id`, `cohort_start_date`, `cohort_end_date`, and `subject_id`. For CDM v5, the `cohort_concept_id` field must be called `cohort_definition_id`.

For our example study, we have created a file called *vignette.sql* with the following contents:

```

/*****
File vignette.sql
*****/

IF OBJECT_ID('@cohortDatabaseSchema.@outcomeTable', 'U') IS NOT NULL
  DROP TABLE @cohortDatabaseSchema.@outcomeTable;

SELECT 1 AS cohort_concept_id,
       condition_start_date AS cohort_start_date,
       condition_end_date AS cohort_end_date,
       condition_occurrence.person_id AS subject_id
INTO @cohortDatabaseSchema.@outcomeTable
FROM @cdmDatabaseSchema.condition_occurrence
INNER JOIN @cdmDatabaseSchema.visit_occurrence
  ON condition_occurrence.visit_occurrence_id = visit_occurrence.visit_occurrence_id
WHERE condition_concept_id IN (
  SELECT descendant_concept_id
  FROM @cdmDatabaseSchema.concept_ancestor
  WHERE ancestor_concept_id = 192671 -- GI - Gastrointestinal haemorrhage
)
AND visit_occurrence.place_of_service_concept_id IN (9201, 9203);

```

This is parameterized SQL which can be used by the `SqlRender` package. We use parameterized SQL so we do not have to pre-specify the names of the CDM and cohort schemas. That way, if we want to run the SQL on a different schema, we only need to change the parameter values; we do not have to change the SQL code. By also making use of translation functionality in `SqlRender`, we can make sure the SQL code can be run in many different environments.

```

library(SqlRender)
sql <- readSql("vignette.sql")
sql <- renderSql(sql,
  cdmDatabaseSchema = cdmDatabaseSchema,
  cohortDatabaseSchema = cohortDatabaseSchema,
  outcomeTable = "my_outcomes")$sql
sql <- translateSql(sql, targetDialect = connectionDetails$dbms)$sql

connection <- connect(connectionDetails)
executeSql(connection, sql)

```

In this code, we first read the SQL from the file into memory. In the next line, we replace the three parameter names with the actual values. We then translate the SQL into the dialect appropriate for the DBMS we already specified in the `connectionDetails`. Next, we connect to the server, and submit the rendered and translated SQL.

If all went well, we now have a table with the outcome of interest. We can see how many events:

```
sql <- paste("SELECT cohort_concept_id, COUNT(*) AS count",
             "FROM @cohortDatabaseSchema.@outcomeTable",
             "GROUP BY cohort_concept_id")
sql <- renderSql(sql,
                 cohortDatabaseSchema = cohortDatabaseSchema,
                 outcomeTable = "my_outcomes")$sql
sql <- translateSql(sql, targetDialect = connectionDetails$dbms)$sql
querySql(connection, sql)

#>   cohort_concept_id  count
#> 1                   1 422274
```

4.3 Extracting the data from the server

Now we can tell `SelfControlledCaseSeries` to extract all necessary data for our analysis:

```
diclofenac <- 1124300

sccsData <- getDbScCsData(connectionDetails = connectionDetails,
                         cdmDatabaseSchema = cdmDatabaseSchema,
                         oracleTempSchema = oracleTempSchema,
                         outcomeDatabaseSchema = cohortDatabaseSchema,
                         outcomeTable = outcomeTable,
                         outcomeIds = 1,
                         exposureDatabaseSchema = cdmDatabaseSchema,
                         exposureTable = "drug_era",
                         exposureIds = diclofenac,
                         cdmVersion = cdmVersion)

sccsData

#> SCCS data object
#>
#> Exposure concept ID(s): 1124300
#> Outcome concept ID(s): 1
```

There are many parameters, but they are all documented in the `SelfControlledCaseSeries` manual. In short, we are pointing the function to the table created earlier and indicating which concept ID in that table identifies the outcome. Note that it is possible to fetch the data for multiple outcomes at once. We further point the function to the `drug_era` table, and specify the concept ID of our exposure of interest: diclofenac. Again, note that it is also possible to fetch data for multiple drugs at once. In fact, when we do not specify any exposure IDs the function will retrieve the data for all the drugs found in the `drug_era` table.

All data about the patients, outcomes and exposures are extracted from the server and stored in the `sccsData` object. This object uses the package `ff` to store information in a way that ensures R does not run out of memory, even when the data are large.

We can use the generic `summary()` function to view some more information of the data we extracted:

```
summary(sccsData)
```

```
#> sccsData object summary
#>
#> Exposure concept ID(s): 1124300
#> Outcome concept ID(s): 1
#>
#> Cases: 211179
#>
#> Outcome counts:
#>   Event count Case count
#> 1      422274      211179
#>
#> Covariates:
#> Number of covariates: 1
#> Number of covariate eras: 18683
```

4.3.1 Saving the data to file

Creating the `sccsData` file can take considerable computing time, and it is probably a good idea to save it for future sessions. Because `sccsData` uses `ff`, we cannot use R's regular save function. Instead, we'll have to use the `saveSccsData()` function:

```
saveSccsData(sccsData, "diclofenacAndGiBleed")
```

We can use the `loadSccsData()` function to load the data in a future session.

4.4 Using a simple model

Next, we can use the data to specify a simple model which we can fit:

```
covarDiclofenac = createCovariateSettings(label = "Exposure of interest",
                                          includeCovariateIds = diclofenac,
                                          start = 0,
                                          end = 0,
                                          addExposedDaysToEnd = TRUE)

sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = covarDiclofenac)

model <- fitSccsModel(sccsEraData)
```

In this example, we use the `createCovariateSettings` to define a single covariate: exposure to diclofenac. We specify that the risk window is from start of exposure to the end by setting `start` and `end` to 0, and requiring that the length of exposure is added to the end date.

We then use the covariate definition in the `createSccsEraData`, and also specify that the first 180 days of observation of every person, the so-called 'naive period', will be excluded from the analysis. Note that data in the naive period will be used to determine exposure status at the start of follow-up (after the end of the naive period). We also specify we will use all occurrences of the outcome, not just the first one per person.

The `fitSccsModel` function is used to fit the model. We can inspect the resulting model:

```
summary(model)
```

```
#> sccsModel object summary
#>
#> Outcome ID: 1
#>
#> Outcome count:
#>   Event count Case count
#> 1      318673    154475
#>
#> Estimates:
#>
#>               Name      ID Estimate  lower .95  upper .95
#> Exposure of interest: Diclofenac 1000      1.367      1.292      1.446
#>   logRr  seLogRr
#> 0.3127 0.02846
```

This tells us what the estimated relative risk (the incidence rate ratio) is during exposure to diclofenac compared to non-exposed time. Note that we lost some cases due to imposing the 180 day naive period.

4.5 Adding a pre-exposure window

The fact that NSAIDs like diclofenac can cause GI bleeds is well known to doctors, and this knowledge affects prescribing behavior. For example, a patient who has just had a GI bleed is not likely to be prescribed diclofenac. This may lead to underestimation of the rate during unexposed time, because the unexposed time includes time just prior to exposure where observing of the outcome is unlikely because of this behavior. One solution to this problem that is often used is to introduce a separate ‘risk window’ just prior to exposure, to separate it from the remaining unexposed time. We can add such a ‘pre-exposure window’ to our analysis:

```
covarPreDiclofenac = createCovariateSettings(label = "Pre-exposure",
                                             includeCovariateIds = diclofenac,
                                             start = -60,
                                             end = -1)

sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = list(covarDiclofenac,
                                                         covarPreDiclofenac))

model <- fitSccsModel(sccsEraData)
```

Here we created a new covariate definition in addition to the first one. We define the risk window to start 60 days prior to exposure, and end on the day just prior to exposure. We combine the two covariate settings in a list for the `createSccsEraData` function. Again, we can take a look at the results:

```
summary(model)
```

```
#> sccsModel object summary
#>
```

```

#> Outcome ID: 1
#>
#> Outcome count:
#>   Event count Case count
#> 1      318673    154475
#>
#> Estimates:
#>
#>           Name      ID Estimate  lower .95  upper .95
#> Exposure of interest: Diclofenac 1000    1.3570    1.2826    1.4348
#>   Pre-exposure: Diclofenac 1001    0.9179    0.8627    0.9756
#>      logRr  seLogRr
#>    0.30528 0.02843
#>   -0.08566 0.03110

```

Here we indeed see a lower relative risk in the time preceding the exposure, indicating the outcome might be a contra-indication for the drug of interest.

4.6 Splitting risk windows

Often we will want to split the risk windows into smaller parts and compute estimates for each part. This can give us insight into the temporal distribution of the risk. We can add this to the model:

```

covarDiclofenacSplit = createCovariateSettings(label = "Exposure of interest",
                                              includeCovariateIds = diclofenac,
                                              start = 0,
                                              end = 0,
                                              addExposedDaysToEnd = TRUE,
                                              splitPoints = c(7,14))

covarPreDiclofenacSplit = createCovariateSettings(label = "Pre-exposure",
                                                  includeCovariateIds = diclofenac,
                                                  start = -60,
                                                  end = -1,
                                                  splitPoints = c(-30))

sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = list(covarDiclofenacSplit,
                                                         covarPreDiclofenacSplit))

```

Here we've redefined our covariate definitions: We kept the same start and end dates, but enforced split points for the main exposure windows at 7 and 14 days. For the pre-exposure window we divided the window into two, at day 30 before the exposure start. Note that the split point dates indicate the end date of the preceding part, so the exposure is now split into day 0 to (and including) day 7, day 8 to (and including) day 14, and day 15 until the end of exposure. The results are:

```
summary(model)
```

```

#> sccsModel object summary
#>
#> Outcome ID: 1

```

```

#>
#> Outcome count:
#>   Event count Case count
#> 1      318673      154475
#>
#> Estimates:
#>
#>      Name      ID Estimate lower .95
#> Exposure of interest: Diclofenac, day 0-7 1000      1.390      1.2194
#> Exposure of interest: Diclofenac, day 8-14 1001      1.429      1.2327
#> Exposure of interest: Diclofenac, day 15- 1002      1.329      1.2443
#>   Pre-exposure: Diclofenac, day -60--30 1003      1.010      0.9323
#>   Pre-exposure: Diclofenac, day -29- 1004      0.830      0.7583
#> upper .95      logRr seLogRr
#>      1.5755      0.32913 0.06401
#>      1.6449      0.35676 0.07191
#>      1.4178      0.28423 0.03309
#>      1.0923      0.01005 0.03994
#>      0.9063     -0.18634 0.04487

```

We see that the risk for the three exposure windows is more or less the same, suggesting a constant risk. We also see that the period 60 to 30 days prior to exposure does not seem to show a decreased risk, suggesting the effect of the contra-indication does not extend more than 30 days before the exposure.

4.7 Including age and seasonality

Often both the rate of exposure and the outcome change with age, and can even depend on the season. This may lead to confounding and may bias our estimates. To correct for this we can include age and/or season into the model.

For computational reasons we assume the effect of both age and season are constant within each calendar month. We assume that the rate from one month to the next can be different, but we also assume that subsequent months have somewhat similar rates. This is implemented by using cubic spline functions.

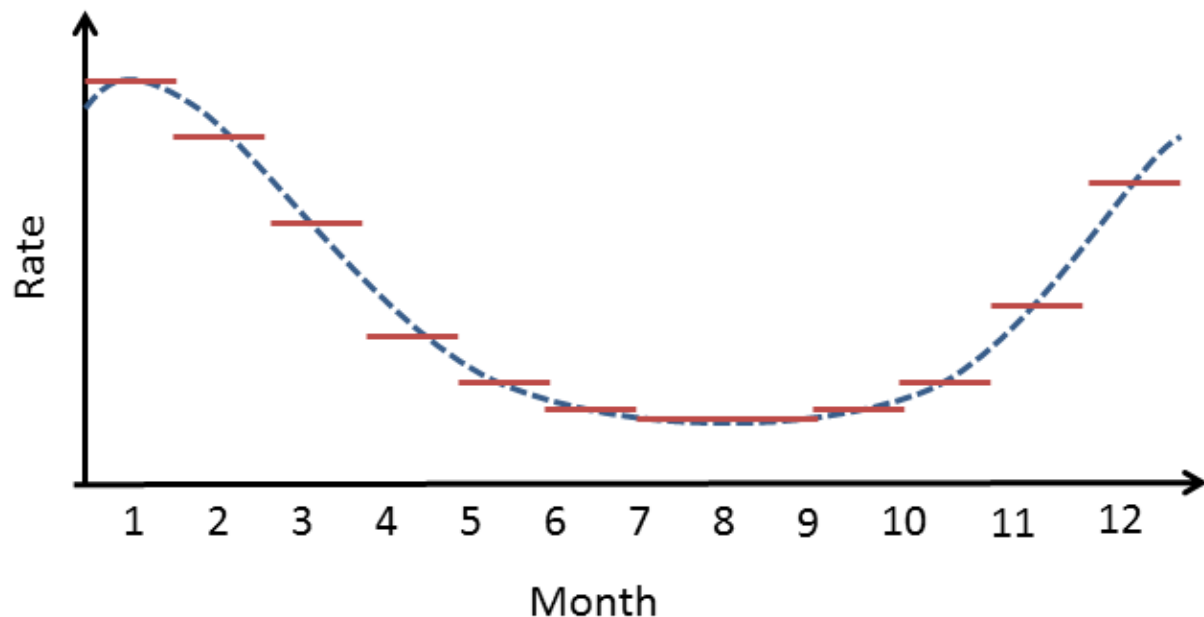


Figure 1. Example of how a spline is used for seasonality: within a month, the risk attributable to seasonality

is assumed to be constant, but from month to month the risks are assumed to follow a cyclic cubic spline.

Note that the by default all people that have the outcome will be used to estimate the effect of age and seasonality on the outcome, so not just the people exposed to the drug of interest. We can add age and seasonality like this:

```
ageSettings <- createAgeSettings(includeAge = TRUE,
                                ageKnots = 5)

seasonalitySettings <- createSeasonalitySettings(includeSeasonality = TRUE,
                                                  seasonKnots = 5)

sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = list(covarDiclofenacSplit,
                                                         covarPreDiclofenacSplit),
                                ageSettings = ageSettings,
                                seasonalitySettings = seasonalitySettings)

model <- fitSccsModel(sccsEraData)
```

Again, we can inspect the model:

```
summary(model)
```

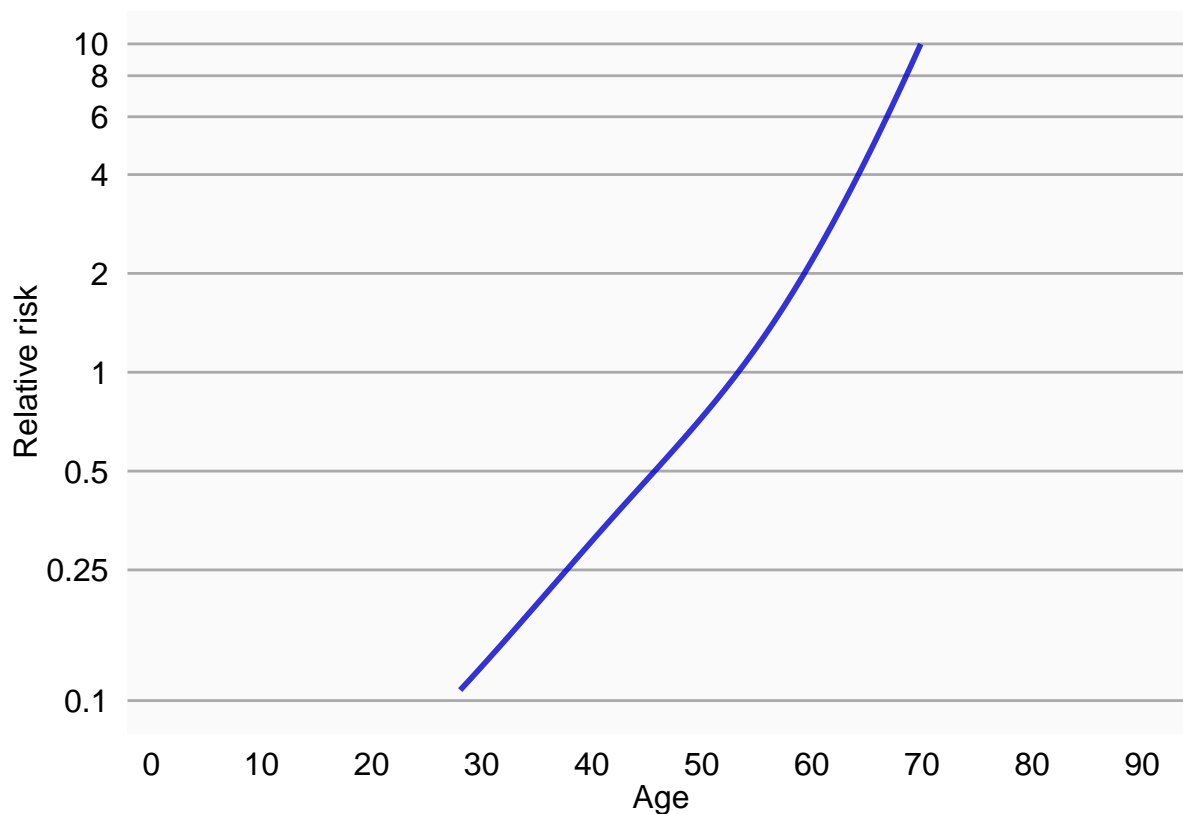
```
#> sccsModel object summary
#>
#> Outcome ID: 1
#>
#> Outcome count:
#>   Event count Case count
#> 1      318673    154475
#>
#> Estimates:
#>
#>               Name      ID  Estimate  lower .95
#>           Age spline component 1    100  6.573e+00  5.080e+00
#>           Age spline component 2    101  4.007e+01  3.205e+01
#>           Age spline component 3    102  4.486e+02  3.489e+02
#>           Age spline component 4    103  4.479e+03  3.458e+03
#>           Age spline component 5    104  2.658e+04  2.057e+04
#>   Seasonality spline component 1    200  1.013e+00  9.835e-01
#>   Seasonality spline component 2    201  1.144e+00  1.126e+00
#>   Seasonality spline component 3    202  9.775e-01  9.465e-01
#> Exposure of interest: Diclofenac, day 0-7 1000  1.381e+00  1.212e+00
#> Exposure of interest: Diclofenac, day 8-14 1001  1.420e+00  1.225e+00
#> Exposure of interest: Diclofenac, day 15- 1002  1.315e+00  1.232e+00
#>   Pre-exposure: Diclofenac, day -60--30 1003  1.013e+00  9.349e-01
#>   Pre-exposure: Diclofenac, day -29- 1004  8.273e-01  7.558e-01
#> upper .95      logRr      seLogRr
#> 8.509e+00    1.88300  0.131697
#> 5.012e+01    3.69069  0.114134
#> 5.770e+02    6.10604  0.128429
#> 5.805e+03    8.40726  0.132266
```

```
#> 3.436e+04 10.18799 0.130955
#> 1.044e+00 0.01339 0.015330
#> 1.162e+00 0.13441 0.008201
#> 1.010e+00 -0.02274 0.016439
#> 1.566e+00 0.32284 0.064075
#> 1.635e+00 0.35048 0.071931
#> 1.403e+00 0.27402 0.033119
#> 1.096e+00 0.01289 0.039967
#> 9.033e-01 -0.18963 0.044866
```

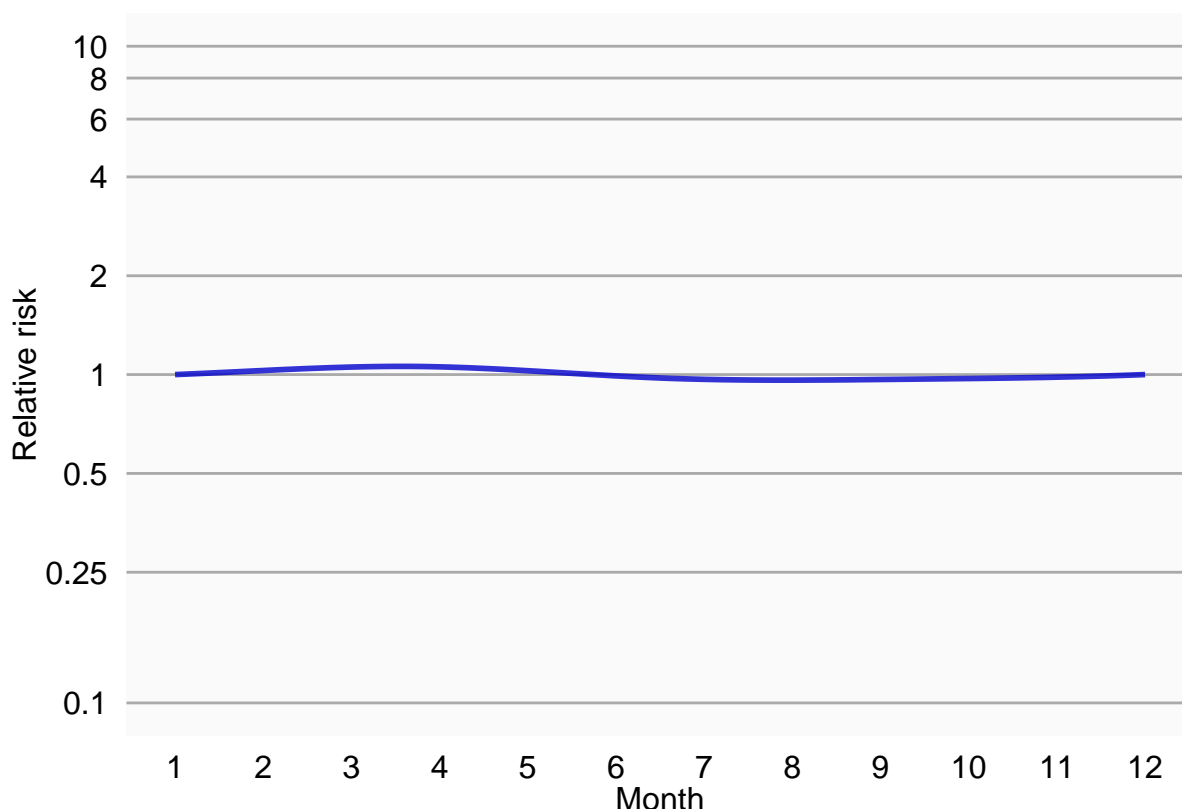
We see that our estimates for exposed and pre-exposure time have not changes much. We can plot the spline curves for age and season to learn more:

```
plotAgeEffect(model)
```

```
#> Warning: Removed 52 rows containing missing values (geom_path).
```



```
plotSeasonality(model)
```



We see a strong effect for age on the outcome, but this effect is spread out over many years and so it less likely to affect the estimates for any individual, since most people are only observed for a few years in the database. We do not see a strong effect for season.

4.8 Considering event-dependent observation time

The SCCS method requires that observation periods are independent of outcome times. This requirement is violated when outcomes increase the mortality rate, since censoring of the observation periods is then event-dependent. A modification to the SCCS has been proposed that attempts to correct for this. First, several models are fitted to estimate the amount and shape of the event-dependent censoring, and the best fitting model is selected. Next, this model is used to reweigh various parts of the observation time. This approach is also implemented in this package, and can be turned on using the `eventDependentObservation` argument of the `createSccsEraData` function:

```
sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = list(covarDiclofenacSplit,
                                                         covarPreDiclofenacSplit),
                                ageSettings = ageSettings,
                                seasonalitySettings = seasonalitySettings,
                                eventDependentObservation = TRUE)

model <- fitSccsModel(sccsEraData)
```

Again, we can inspect the model:

```
summary(model)
```

```
#> sccsModel object summary
#>
#> Outcome ID: 1
#>
#> Outcome count:
#>   Event count Case count
#> 1      318673    154475
#>
#> Estimates:
#>
#>               Name      ID  Estimate  lower .95
#>      Age spline component 1    100  5.451e-03  4.235e-03
#>      Age spline component 2    101  1.121e-04  9.012e-05
#>      Age spline component 3    102  5.857e-06  4.585e-06
#>      Age spline component 4    103  2.377e-06  1.849e-06
#>      Age spline component 5    104  5.050e-07  3.940e-07
#>      Seasonality spline component 1    200  9.821e-01  9.533e-01
#>      Seasonality spline component 2    201  1.149e+00  1.131e+00
#>      Seasonality spline component 3    202  9.879e-01  9.568e-01
#>      Exposure of interest: Diclofenac, day 0-7    1000  1.408e+00  1.236e+00
#>      Exposure of interest: Diclofenac, day 8-14    1001  1.443e+00  1.245e+00
#>      Exposure of interest: Diclofenac, day 15-    1002  1.342e+00  1.257e+00
#>      Pre-exposure: Diclofenac, day -60--30    1003  1.053e+00  9.722e-01
#>      Pre-exposure: Diclofenac, day -29-    1004  8.512e-01  7.775e-01
#> upper .95      logRr      seLogRr
#> 7.015e-03    -5.21199  0.128693
#> 1.393e-04    -9.09650  0.111104
#> 7.477e-06   -12.04785  0.124567
#> 3.056e-06   -12.94948  0.128130
#> 6.473e-07   -14.49868  0.126632
#> 1.012e+00    -0.01806  0.015174
#> 1.168e+00     0.13891  0.008218
#> 1.020e+00    -0.01214  0.016349
#> 1.597e+00     0.34248  0.064118
#> 1.661e+00     0.36680  0.071877
#> 1.433e+00     0.29435  0.033223
#> 1.139e+00     0.05201  0.039932
#> 9.296e-01    -0.16113  0.044976
```

5 Studies with more than one drug

Although we are usually interested in the effect of a single drug or drug class, it could be beneficial to add exposure to other drugs to the analysis if we believe those drugs represent time-varying confounders that we wish to correct for.

5.1 Adding a class of drugs

For example, oftentimes diclofenac is co-prescribed with proton-pump inhibitors (PPIs) to mitigate the risk of GI bleeding. We would like our estimate to represent just the effect of the diclofenac, so we need to keep the effect of the PPIs separate. First we have to retrieve the information on PPI exposure from the database:

```

diclofenac <- 1124300
ppis <- c(911735, 929887, 923645, 904453, 948078, 19039926)

sccsData <- getDbSccsData(connectionDetails = connectionDetails,
                          cdmDatabaseSchema = cdmDatabaseSchema,
                          oracleTempSchema = oracleTempSchema,
                          outcomeDatabaseSchema = cohortDatabaseSchema,
                          outcomeTable = outcomeTable,
                          outcomeIds = 1,
                          exposureDatabaseSchema = cdmDatabaseSchema,
                          exposureTable = "drug_era",
                          exposureIds = c(diclofenac, ppis),
                          cdmVersion = cdmVersion)

sccsData

#> SCCS data object
#>
#> Exposure concept ID(s): 1124300,911735,929887,923645,904453,948078,19039926
#> Outcome concept ID(s): 1

```

Once retrieved, we can use the data to build and fit our model:

```

covarPpis = createCovariateSettings(label = "PPIs",
                                   includeCovariateIds = ppis,
                                   stratifyByID = FALSE,
                                   start = 1,
                                   end = 0,
                                   addExposedDaysToEnd = TRUE)

sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = list(covarDiclofenacSplit,
                                                         covarPreDiclofenacSplit,
                                                         covarPpis),
                                ageSettings = ageSettings,
                                seasonalitySettings = seasonalitySettings,
                                eventDependentObservation = TRUE)

model <- fitSccsModel(sccsEraData)

```

Here, we added a new covariate based on the list of concept IDs for the various PPIs. In this example we set **stratifyByID** to **FALSE**, meaning that we will estimate a single incidence rate ratio for all PPIs, so one estimate for the entire class of drugs. Note that duplicates will be removed: if a person is exposed to two PPIs on the same day, this will be counted only once when fitting the model. Furthermore, we have set the **start** day to 1 instead of 0. The reason for this is that PPIs will also be used to treat GI bleeds, and are likely to be prescribed on the same day as the event. If we would include day 0, the risk of the outcome would be attributed to the PPI used for treatment, not the other factors that caused the GI bleed such as any exposure to our drug of interest. Again, we can inspect the model:

```
summary(model)
```

```
#> sccsModel object summary
```

```

#>
#> Outcome ID: 1
#>
#> Outcome count:
#>   Event count Case count
#> 1      318673    154475
#>
#> Estimates:
#>
#>           Name      ID  Estimate  lower .95
#>           Age spline component 1    100  5.590e-03  4.344e-03
#>           Age spline component 2    101  1.221e-04  9.820e-05
#>           Age spline component 3    102  6.482e-06  5.076e-06
#>           Age spline component 4    103  2.661e-06  2.069e-06
#>           Age spline component 5    104  5.666e-07  4.419e-07
#>           Seasonality spline component 1    200  9.807e-01  9.518e-01
#>           Seasonality spline component 2    201  1.150e+00  1.131e+00
#>           Seasonality spline component 3    202  9.875e-01  9.563e-01
#> Exposure of interest: Diclofenac, day 0-7    1000  1.420e+00  1.246e+00
#> Exposure of interest: Diclofenac, day 8-14    1001  1.456e+00  1.256e+00
#> Exposure of interest: Diclofenac, day 15-    1002  1.358e+00  1.271e+00
#> Pre-exposure: Diclofenac, day -60--30    1003  1.057e+00  9.758e-01
#> Pre-exposure: Diclofenac, day -29-    1004  8.558e-01  7.815e-01
#> PPIs    1005  8.639e-01  8.507e-01
#> upper .95      logRr      seLogRr
#> 7.194e-03    -5.18669  0.128689
#> 1.519e-04    -9.01047  0.111193
#> 8.277e-06   -11.94645  0.124686
#> 3.422e-06   -12.83677  0.128271
#> 7.264e-07   -14.38366  0.126765
#> 1.010e+00    -0.01948  0.015260
#> 1.168e+00     0.13953  0.008217
#> 1.020e+00    -0.01260  0.016355
#> 1.610e+00     0.35072  0.064123
#> 1.677e+00     0.37556  0.072038
#> 1.449e+00     0.30576  0.033185
#> 1.144e+00     0.05582  0.040009
#> 9.346e-01    -0.15576  0.044986
#> 8.774e-01    -0.14625  0.007883

```

We do see a decrease in risk when people are exposed to PPIs.

5.2 Adding all drugs

Another approach could be to add all drugs into the model. Again, the first step is to get all the relevant data from the database:

```

sccsData <- getDbScCsData(connectionDetails = connectionDetails,
                           cdmDatabaseSchema = cdmDatabaseSchema,
                           oracleTempSchema = oracleTempSchema,
                           outcomeDatabaseSchema = cohortDatabaseSchema,
                           outcomeTable = outcomeTable,
                           outcomeIds = 1,
                           exposureDatabaseSchema = cdmDatabaseSchema,

```

```

exposureTable = "drug_era",
exposureIds = c(),
cdmVersion = cdmVersion)

```

Note that the `exposureIds` argument is left empty. This will cause data for all concepts in the exposure table to be retrieved. Next, we simply create a new set of covariates, and fit the model:

```

covarAllDrugs = createCovariateSettings(label = "All other exposures",
                                       excludeCovariateIds = diclofenac,
                                       stratifyByID = TRUE,
                                       start = 1,
                                       end = 0,
                                       addExposedDaysToEnd = TRUE,
                                       allowRegularization = TRUE)

sccsEraData <- createSccsEraData(sccsData,
                                naivePeriod = 180,
                                firstOutcomeOnly = FALSE,
                                covariateSettings = list(covarDiclofenacSplit,
                                                         covarPreDiclofenacSplit,
                                                         covarAllDrugs),
                                ageSettings = ageSettings,
                                seasonalitySettings = seasonalitySettings,
                                eventDependentObservation = TRUE)

model <- fitSccsModel(sccsEraData)

```

The first thing to note is that we have defined the new covariates to be all drugs except diclofenac by not specifying the `includeCovariateIds` and setting the `excludeCovariateIds` to the concept ID of diclofenac. Furthermore, we have specified that `stratifyByID` is `TRUE`, meaning an estimate will be produced for each drug.

We have set `allowRegularization` to `TRUE`, meaning we will use regularization for all estimates in this new covariate set. Regularization means we will impose a prior distribution on the effect size, effectually penalizing large estimates. This helps fit the model, for example when some drugs are rare, and when drugs are almost often prescribed together and their individual effects are difficult to untangle.

Because there are now so many estimates, we will not use the `summary()` function but instead export all estimates to a data frame using `getModel()`:

```

estimates <- getModel(model)
estimates[estimates$originalCovariateId == diclofenac, ]

```

```

#>           name    id estimate    lb95Ci
#> 9  Exposure of interest: Diclofenac, day 0-7 1000 1.3106430 1.1499519
#> 10 Exposure of interest: Diclofenac, day 8-14 1001 1.3718646 1.1831907
#> 11 Exposure of interest: Diclofenac, day 15- 1002 1.2795529 1.1967493
#> 12   Pre-exposure: Diclofenac, day -60--30 1003 1.0329653 0.9530873
#> 13   Pre-exposure: Diclofenac, day -29- 1004 0.8189428 0.7477657
#>      ub95Ci    logRr    seLogRr originalCovariateId
#> 9  1.4867372 0.2705179 0.06432060             1124300
#> 10 1.5800021 0.3161708 0.07207039             1124300
#> 11 1.3667779 0.2465107 0.03364618             1124300

```

```
#> 12 1.1175772 0.0324336 0.04016887 1124300
#> 13 0.8947282 -0.1997411 0.04515684 1124300
#> originalCovariateName
#> 9 Diclofenac
#> 10 Diclofenac
#> 11 Diclofenac
#> 12 Diclofenac
#> 13 Diclofenac
```

Here we see that despite the extensive adjustments that are made in the model, the effect estimates for diclofenac have remained nearly the same.

In case we're interested, we can also look at the effect sizes for the PPIs:

```
estimates[estimates$originalCovariateId %in% ppis, ]
```

```
#> name id estimate lb95Ci ub95Ci
#> 88 Other exposures: Esomeprazole 1147 0.6790260 NA NA
#> 98 Other exposures: rabeprazole 1171 0.8570214 NA NA
#> 115 Other exposures: Omeprazole 1203 0.7675711 NA NA
#> 121 Other exposures: lansoprazole 1219 0.8275570 NA NA
#> 136 Other exposures: pantoprazole 1257 0.7658544 NA NA
#> 432 Other exposures: dextansoprazole 1860 0.7085762 NA NA
#> logRr seLogRr originalCovariateId originalCovariateName
#> 88 -0.3870959 NA 904453 Esomeprazole
#> 98 -0.1542923 NA 911735 rabeprazole
#> 115 -0.2645242 NA 923645 Omeprazole
#> 121 -0.1892773 NA 929887 lansoprazole
#> 136 -0.2667631 NA 948078 pantoprazole
#> 432 -0.3444977 NA 19039926 dextansoprazole
```

Note that because we used regularization, we are not able to compute the confidence intervals for these estimates. We do again see that PPIs all have relative risks lower than 1 as we would expect.

6 Acknowledgments

Considerable work has been dedicated to provide the `SelfControlledCaseSeries` package.

```
citation("SelfControlledCaseSeries")
```

```
#>
#> To cite package 'SelfControlledCaseSeries' in publications use:
#>
#> Martijn Schuemie, Patrick Ryan, Marc Suchard and Trevor Shaddox
#> (2015). SelfControlledCaseSeries: Self-Controlled Case Series. R
#> package version 0.0.2.
#>
#> A BibTeX entry for LaTeX users is
#>
#> @Manual{,
#> title = {SelfControlledCaseSeries: Self-Controlled Case Series},
```



```

#>   author = {Martijn Schuemie and Patrick Ryan and Marc Suchard and Trevor Shaddox},
#>   year = {2015},
#>   note = {R package version 0.0.2},
#> }
#>
#> ATTENTION: This citation information has been auto-generated from
#> the package DESCRIPTION file and may need manual editing, see
#> 'help("citation")'.

```

Furthermore, `SelfControlledCaseSeries` makes extensive use of the `Cyclops` package.

```
citation("Cyclops")
```

```

#>
#> To cite Cyclops in publications use:
#>
#> Suchard MA, Simpson SE, Zorych I, Ryan P and Madigan D (2013).
#> "Massive parallelization of serial inference algorithms for
#> complex generalized linear models." _ACM Transactions on Modeling
#> and Computer Simulation_, *23*, pp. 10. <URL:
#> http://dl.acm.org/citation.cfm?id=2414791>.
#>
#> A BibTeX entry for LaTeX users is
#>
#> @Article{,
#>   author = {M. A. Suchard and S. E. Simpson and I. Zorych and P. Ryan and D. Madigan},
#>   title = {Massive parallelization of serial inference algorithms for complex generalized linear m
#>   journal = {ACM Transactions on Modeling and Computer Simulation},
#>   volume = {23},
#>   pages = {10},
#>   year = {2013},
#>   url = {http://dl.acm.org/citation.cfm?id=2414791},
#> }

```

Part of the code (related to event-dependent observation periods) is based on the `SCCS` package by Yonas Ghebremichael-Weldeselassie, Heather Whitaker, and Paddy Farrington.

This work is supported in part through the National Science Foundation grant IIS 1251151.