

# Package: aciccomp2016 (via r-universe)

May 28, 2026

**Version** 0.1-0

**Date** 2017-05-20

**Title** Atlantic Causal Inference Conference Competition 2016 Simulation

**Depends** R (>= 3.2-2)

**Imports** stats, methods, utils

**Suggests** testthat

**Description** Generate simulation data.

**License** GPL (>= 2)

**NeedsCompilation** no

**Biarch** yes

**LazyData** yes

**Repository** <https://mrcieu.r-universe.dev>

**Date/Publication** 2020-07-08 20:56:10 UTC

**RemoteUrl** <https://github.com/vdorie/aciccomp>

**RemoteRef** HEAD

**RemoteSha** 282d26659b2d3d6fd060dde6d32feeb8f1e8ab5a

**RemoteSubdir** 2016

## Contents

constants_2016 . . . . .	2
dgp_2016 . . . . .	5
input_2016 . . . . .	7
parameters_2016 . . . . .	9

<b>Index</b>	<b>11</b>
--------------	-----------

---

 constants\_2016

*Constants Used in DGP for ACIC Competition 2016*


---

**Description**

Returns or sets elements of a named list containing all of the constants required to run the data generating processes for the 2016 ACIC Competition.

**Usage**

```
constants_2016(...)
```

**Arguments**

... Options from the list below.

**Details**

Returns default values or sets them, as appropriate. Minimal error checking is performed.

**Value**

RSP_INPUT_SCALE	Scaling factor applied to covariates before evaluating the response function.
RSP_OUTPUT_SHAPE_1	The first shape parameter in a beta-prime used to generate the output scale of the response function.
RSP_OUTPUT_RATE	The inverse scale parameter in a beta-prime used to generate the output scale of the response function.
RSP_OUTPUT_SHAPE_2	The second shape parameter in a beta-prime used to generate the output scale of the response function.
TRT_INPUT_SCALE	Scaling factor applied to covariates before evaluating the treatment assignment function.
TRT_OUTPUT_SCALE	Scaling factor applied to result of the treatment assignment function.
TRT_BIAS_SCALE	Approximate scale for treatment biasing functions when overlap parameter is not "full".
RSP_SIGMA_Y	Scale of noise added to response.
BF_CONSTANT_SCALE	Scale of constant base function parameter.
BF_LINEAR_SCALE	Scale of linear base function parameter.

BF_QUADRATIC_SHAPE_1	First shape parameter used to generate quadratic base function root parameter.
BF_QUADRATIC_SHAPE_2	Second shape parameter used to generate quadratic base function root parameter.
BF_QUADRATIC_RATE	Rate parameter used to generate quadratic base function root parameter.
BF_QUADRATIC_SCALE	Scale of quadratic base function parameter.
BF_CUBIC_SHAPE	Shape parameter used to generate cubic base function root parameters.
BF_CUBIC_RATE	Rate parameter used to generate cubic base function root parameters.
BF_CUBIC_SCALE	Scale of cubic base function parameter.
BF_CONTINUOUS_SCALE	Scale parameter shared by continuous base functions.
BF_STEP_SHAPE	Shape parameter for step base functions.
BF_STEP_CONSTANT_SCALE	Scale of step-wise constant base function parameter.
BF_STEP_LINEAR_SCALE	Scale of piece-wise linear base function parameter.
BF_SIGMOID_SHAPE_1	First shape parameter used to generate sigmoid base function parameters.
BF_SIGMOID_RATE_1	First rate parameter used to generate sigmoid base function parameters.
BF_SIGMOID_SHAPE_2	Second shape parameter used to generate sigmoid base function parameters.
BF_SIGMOID_RATE_2	Second rate parameter used to generate sigmoid base function parameters.
BF_QUANTILE_SHAPE_1	First shape parameter used to generate quantile base function cutoff.
BF_QUANTILE_SHAPE_2	Second shape parameter used to generate quantile base function cutoff.
BF_TWEAK_SIGN_PROB	Probability of changing sign when copy/modifying base function.
BF_TWEAK_NORMAL_SCALE	Scale of normal noise added to unconstrained base function parameters when copy/modifying.
BF_TWEAK_GAMMA_SHAPE	Shape parameter of positive noise added to constrained base function parameters when copying/modifying.
BF_TWEAK_GAMMA_RATE	Rate parameter of positive noise added to constrained base function parameters when copying/modifying.
TRT_BF_DF	Base function degrees of freedom when generating treatment assignment mechanism.
RSP_BF_DF	Base function degrees of freedom when generating response surface.

TRT_LINEAR_SCALE_SHAPE_1	First scale parameter used to generate overall scale of treatment assignment mechanism.
TRT_LINEAR_SCALE_SHAPE_2	Second scale parameter used to generate overall scale of treatment assignment mechanism.
TRT_LINEAR_SCALE_RATE	Rate parameter used to generate overall scale of treatment assignment mechanism.
RSP_EXP_SCALE_SHAPE	Shape parameter used when generating scale factor for exponential functions.
RSP_EXP_SCALE_RATE	Rate parameter used when generating scale factor for exponential functions.
RSP_EXP_WEIGHT_SHAPE	Shape parameter used when generating relative weight factor for exponential functions.
RSP_EXP_WEIGHT_RATE	Rate parameter used when generating relative weight factor for exponential functions.
RSP_TE_MEAN	Expected value for population average treatment effect.
RSP_TE_SCALE	Scale factor for population average treatment effect.
RSP_TE_DF	Degrees of freedom for population average treatment effect.
SPARSE_COVARIATE_WEIGHT	Weight of inclusion for sparse, discrete covariates.
CONTINUOUS_COVARIATE_WEIGHT	Weight of inclusion for continuous covariates.
DEFAULT_COVARIATE_WEIGHT	Default weight of inclusion for covariates.
TRT_BASELINE_SHIFT	Function used to derive a scale when generating a baseline treatment probability from <code>root.trt</code>
BASE_FUNCTION_DIST_LIN	Base function distribution containing only linear functions.
BASE_FUNCTION_DIST_POLY	Base function distribution containing linear, quadratic, and cubic functions.
BASE_FUNCTION_DIST_STEP	Base function distribution containing linear, step-wise constant, and piece-wise linear functions.
BASE_FUNCTION_DIST_EXP	Base function distribution containing third order polynomials to be used in exponential functions.
<code>dist.lin</code>	Function distribution for purely linear treatment or response.
<code>dist.int</code>	Function distribution with linear terms and interactions.
<code>dist.pure.poly</code>	Function distribution with quadratic terms and no interactions.
<code>dist.poly</code>	Function distribution with cubic terms and interactions.

dist.step	Function distribution with linear terms, step-wise constant terms, and interactions.
dist.exp	Function distribution with quadratic terms and interactions appropriate for use with exponential link functions.
dist.bias1	Function distribution over treatment assignment biasing functions.
dist.bias2	Function distribution over treatment assignment biasing functions.
dist.hetero.med	Function distribution specifying interaction retention probabilities for medium degrees of treatment effect heterogeneity.
dist.hetero.high	Function distribution specifying interaction retention probabilities for high degrees of treatment effect heterogeneity.

**Author(s)**

Vincent Dorie: <vdorie@gmail.com>.

---

dgp\_2016

*Data Generating Process for the 2016 ACIC Competition*


---

**Description**

Applies the data generating process used in the Atlantic Causal Inference Competition of 2016.

**Usage**

```
dgp_2016(x, parameters, random.seed,
         constants = constants_2016(),
         extraInfo = FALSE)
```

**Arguments**

x	Input data in the form of a data frame, most likely <a href="#">input_2016</a> .
parameters	A named list containing elements in the form of <a href="#">parameters_2016</a> , a row of the same object, or an integer specifying which row of <a href="#">parameters_2016</a> is to be used; see that page for details.
random.seed	A list of arguments to be used in a call to <a href="#">set.seed</a> or an integer between 1 and 100 specifying the random seed associated with an iteration from the competition.
constants	A named list containing elements as returned by <a href="#">constants_2016</a> ; see there for details.
extraInfo	Boolean determining if additional information is to be returned, including the treatment and control response surfaces, the transformed input data, and whether or not a simulation would have been deemed interesting enough to include in the competition.

## Details

Creates a causal inference problem by taking the input  $x$  and using the passed in parameters to generate a treatment assignment mechanism (probability of treatment for each individual), response surface (expected value under treatment and control), and finally observed data. The parameters provide high-level controls to adjust the result for causal inference features that may be of interest, while constants control at a lower level the parameters of generated functions.

**Generalized Additive Functions:** The 2016 competition used a unique set of software that was internally described as “Generalized Additive Functions” (GAFs). A GAF consists of many small functions applied to various features/columns of the input that are added together or interacted with each other. The complete sum may then be passed through a link function to achieve a result in a transformed space. The small functions are randomly derived from a library of functions, so that the general features of the result can vary according to high level parameters.

This package reproduces GAFs as they were used in the 2016 contest without the intention that they be further applied. It may be possible to use `dgp_2016` with different input data and changes to the constants should propagate, however these extensions will not be widely supported.

## Value

A named list containing:

<code>z</code>	Vector of treatment assignments. If <code>extraInfo</code> is FALSE, <code>z</code> contains 0s and 1s. If TRUE, <code>z</code> is a factor with levels <code>ctl</code> and <code>trt</code> .
<code>y</code>	Vector of observed response variables, $y(z)$ .
<code>y.0</code>	Vector of response variables under the control condition, $y(0)$ .
<code>y.1</code>	Vector of response variables under the treatment condition, $y(1)$ .
<code>mu.0</code>	Vector of expected response under the control condition, $E[Y(0)]$ .
<code>mu.1</code>	Vector of expected response under the treatment condition, $E[Y(1)]$ .
<code>e</code>	Vector of propensity scores, $P(Z = 1)$ .
<code>f.z</code>	Optional - the GAF for the treatment assignment mechanism.
<code>f.y</code>	Optional - the GAF for the response surface.
<code>x</code>	Optional - the transformed input passed to <code>f.z</code> and <code>f.y</code> .
<code>valid</code>	Optional - boolean if the simulation would be rejected as "uninteresting".

## Author(s)

Vincent Dorie: <vdorie@gmail.com>.

## References

Dorie V., Hill J., Shalit U., Scott M. and Cervone D. (2017) Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition, preprint arXiv <https://arxiv.org/abs/1707.02641>.

## Examples

```
## Not run:
# to test a method
ate <- matrix(NA, 77, 100)
for (i in seq_len(77)) {
  for (j in seq_len(100)) {
    sim <- dgp_2016(input_2016, i, j)
    df <- input_2016
    df$y <- sim$y
    df$z <- sim$z
    fit <- lm(y ~ ., df)
    ate[i,j] <- coef(fit)["z"]
  }
}

## undocumented features, getting closest approximate linear model
sim <- dgp_2016(input_2016, 1, 1, extraInfo = TRUE)

e <- aciccomp::evaluate(sim$f.z, sim$x)
x.z.approx <- aciccomp::evaluateGeneralizedAdditiveFunctionToDataframe(sim$f.z, sim$x)

x.temp <- sim$x
x.temp$.z <- sim$z
x.y.approx <- aciccomp::evaluateGeneralizedAdditiveFunctionToDataframe(sim$f.y, x.temp)

## End(Not run)
```

---

input\_2016

*Input Data for the 2016 ACIC Competition*

---

## Description

Input data used in the 2016 Atlantic Causal Inference Competition, taken from the Collaborative Perinatal Project.

## Usage

```
input_2016
```

## Format

A data frame consisting of 4802 observations and 58 covariates. The columns have been de-identified from their original source, but correspond to possible confounders, instruments, and uncorrelated variables from a hypothetical twin study on the impact of birthweight on IQ.

**Details**

The variable in the original CPP are:

- mom\_age
- mar\_status
- mom\_cigs\_per\_day
- mom\_years\_smoked
- mom\_height
- mom\_weight\_prior
- mom\_num\_cardio\_cond
- mom\_num\_pulm\_cond
- mom\_num\_hema\_cond
- mom\_num\_endocrine\_cond
- mom\_num\_veneral\_cond
- mom\_num\_urin\_cond
- mom\_num\_gyne\_cond
- mom\_num\_neur\_cond
- mom\_num\_obst\_compl
- mom\_num\_infect\_dis
- mom\_work\_status
- mom\_years\_educ
- family\_income
- housing\_density
- mom\_birth\_place
- consanguinity
- socio\_eco
- mom\_race
- age\_menarche
- dias\_blood\_pres
- mom\_weight\_birth
- dad\_age
- dad\_years\_educ
- num\_premes
- num\_abortions
- num\_prior\_pregs
- num\_stillbirths
- bayley\_mental
- bayley\_motor

- placental\_weight
- cord\_length
- sex
- apgar\_1m\_total
- apgar\_5m\_total
- bottle\_feed\_days
- breast\_feed\_days
- child\_bilirubin
- child\_hematocrit
- child\_hemoglobin
- child\_num\_neur\_abn
- child\_num\_cns\_cond
- child\_num\_muscoskel
- child\_num\_resp\_abn
- child\_num\_cardio\_abn
- child\_num\_liver\_abn
- child\_num\_hemo\_cond
- child\_num\_infect
- child\_num\_synd
- child\_num\_endo\_dis
- child\_num\_proc
- head\_size\_1yr
- gest\_delivery

### Source

Niswander, K. R. and Gordon, M. (1972) The Collaborative Perinatal Study of the National Institute of Neurological Diseases and Stroke: the women and their pregnancies. Philadelphia, PA: W.B. Saunders Company <https://www.archives.gov/research/electronic-records/nih.html>

---

parameters\_2016

*Parameters Data for the 2016 ACIC Competition*

---

### Description

Data set containing the parameters used to generate data for the 2016 Atlantic Causal Inference Conference competition.

### Usage

parameters\_2016

**Format**

A data frame describing 77 scenarios that vary across 6 features.

1. `model.trt` - Function distribution over the treatment assignment mechanism. Can be "linear", "polynomial", or "step".
2. `root.trt` - Baseline probability of receiving treatment.
3. `overlap.trt` - Term that controls the addition of overlap-penalizing terms that forcibly exclude observations from the treatment group by carving out hyper-rectangles of the covariate space and assigning their treatment probability to 0. Can be "full" for complete overlap, "one-term" for adding a single function as described above, or "two-term" for adding two. Two-terms were not used in the competition and is not thoroughly tested.
4. `model.rsp` - Function distribution over the response surface. Can be "linear", "polynomial", "step", or "exponential".
5. `alignment` - A numeric value that determines the degree to which terms from the treatment assignment function appear in response surface function.
6. `te.hetero` - A term that controls the degree of treatment effect heterogeneity. Can be "none" for parallel surfaces, "med" or "high". Higher heterogeneity is achieved by selectively interacting terms from the response surface with a treatment indicator.

**Source**

Original release.

**References**

Dorie V., Hill J., Shalit U., Scott M. and Cervone D. (2017) Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition, preprint arXiv <https://arxiv.org/abs/1707.02641>.

# Index

- \* **causal inference**

- constants\_2016, [2](#)

- dgp\_2016, [5](#)

- \* **datasets**

- input\_2016, [7](#)

- parameters\_2016, [9](#)

- \* **simulation**

- constants\_2016, [2](#)

- dgp\_2016, [5](#)

constants\_2016, [2](#), [5](#)

dgp\_2016, [5](#)

input\_2016, [5](#), [7](#)

parameters\_2016, [5](#), [9](#)

set.seed, [5](#)