

---

# **Bayesian inference of COVID-19**

*Release 0.1.7*

**Jonas Dehning, Johannes Zierenberg, F. Paul Spitzner, Michael W**

**Jun 17, 2020**



# CONTENTS

<b>1</b>	<b>Installation</b>	<b>1</b>
<b>2</b>	<b>First Steps</b>	<b>3</b>
<b>3</b>	<b>Disclaimer</b>	<b>5</b>
<b>4</b>	<b>Model</b>	<b>7</b>
<b>5</b>	<b>Data Retrieval</b>	<b>15</b>
5.1	Utility . . . . .	15
5.2	Johns Hops University . . . . .	16
5.3	Robert Koch Institute . . . . .	18
5.4	Robert Koch Institute situation reports . . . . .	20
5.5	Google . . . . .	21
5.6	Our World in Data . . . . .	22
5.7	Base Retrieval Class . . . . .	23
<b>6</b>	<b>Plotting</b>	<b>25</b>
<b>7</b>	<b>Variables saved in the trace</b>	<b>27</b>
<b>8</b>	<b>Contributing</b>	<b>29</b>
8.1	Beginning . . . . .	29
8.2	Code formatting . . . . .	29
8.3	Testing . . . . .	29
8.4	Documentation . . . . .	30
<b>9</b>	<b>Debugging</b>	<b>31</b>
9.1	General approach for nans/infs during sampling . . . . .	31
9.2	Sampler: MCMC (Nuts) . . . . .	32
9.3	Sampler: Variational Inference . . . . .	33
<b>10</b>	<b>Indices and tables</b>	<b>35</b>
	<b>Bibliography</b>	<b>37</b>
	<b>Python Module Index</b>	<b>39</b>
	<b>Index</b>	<b>41</b>



## INSTALLATION

There exists three different possibilities to run the models:

1. Clone the repository, with the latest release:

```
git clone --branch v0.1.7 https://github.com/Priesemann-Group/covid19_inference
```

2. Install the module via pip

```
pip install git+https://github.com/Priesemann-Group/covid19_inference.git@v0.1.7
```

3. Run the notebooks directly in Google Colab. At the top of the notebooks files there should be a symbol which opens them directly in a Google Colab instance.



## FIRST STEPS

To get started, we recommend to look at one of the currently two example notebooks:

1. **SIR model with one german state** This model is similar to the one discussed in our paper: [Inferring COVID-19 spreading rates and potential change points for case number forecasts](#). The difference is that the delay between infection and report is now lognormal distributed and not fixed.
2. **Hierarchical model of the German states** This builds a hierarchical bayesian model of the states of Germany

We can for example recommend the following articles about bayesian modeling:

As a introduction to Bayesian statistics and the python package (PyMC3) that we use: [https://docs.pymc.io/notebooks/api\\_quickstart.html](https://docs.pymc.io/notebooks/api_quickstart.html)

This is a good post about hierarchical Bayesian models in general: <https://statmodeling.stat.columbia.edu/2014/01/21/everything-need-know-bayesian-statistics-learned-eight-schools/>





## DISCLAIMER

We evaluate the data provided by the John Hopkins University [link](#). We exclude any liability with regard to the quality and accuracy of the data used, and also with regard to the correctness of the statistical analysis. The evaluation of the different growth phases represents solely our personal opinion.

The number of cases reported may be significantly lower than the number of people actually infected. Also, we must point out that week-ends and changes in the test system may lead to fluctuations in reported cases that have no equivalent in actual case numbers.

Certainly, at this stage all statistical predictions are subject to great uncertainty because the general trends of the epidemic are not yet clear. In any case, the statistical trends that we interpret from the data are only suitable for predictions if the measures taken by the government and authorities to contain the pandemic remain in force and are being followed by the population. We must also point out that, even if the statistics indicate that the epidemic is under control, we may at any time see a resurgence of infection figures until the disease is eradicated worldwide.



## MODEL

```
class covid19_inference.model.Cov19Model (new_cases_obs,    data_begin,    fcast_len,
                                         diff_data_sim,    N_population,    name="",
                                         model=None)
```

Model class used to create a covid-19 propagation dynamics model. Parameters below are passed to the constructor. Attributes (Variables) are available after creation and can be accessed from every instance. Some background:

- The simulation starts *diff\_data\_sim* days before the data.
- The data has a certain length, on which the inference is based. This length is given by *new\_cases\_obs*.
- After the inference, a forecast takes of length *fcast\_len* takes place, starting on the day after the last data point in *new\_cases\_obs*.
- In total, traces produced by a model run have the length  $sim\_len = diff\_data\_sim + data\_len + fcast\_len$
- Date ranges include both boundaries. For example, if *data\_begin* is March 1 and *data\_end* is March 3 then *data\_len* will be 3.

### Parameters

- **new\_cases\_obs** (*1 or 2d array*) – If the array is two-dimensional, an hierarchical model will be constructed. First dimension is then time, the second the region/country.
- **data\_begin** (*datetime.datetime*) – Date of the first data point
- **fcast\_len** (*int*) – Number of days the simulations runs longer than the data
- **diff\_data\_sim** (*int*) – Number of days the simulation starts earlier than the data. Should be significantly longer than the delay between infection and report of cases.
- **N\_population** (*number or 1d array*) – Number of inhabitation in region, needed for the S(E)IR model. Is ideally 1 dimensional if new\_cases\_obs is 2 dimensional
- **name** (*string*) – suffix appended to the name of random variables saved in the trace
- **model** – specify a model, if this one should expand another

### Variables

- **new\_cases\_obs** (*1 or 2d array*) – as passed during construction
- **data\_begin** (*datetime.datetime*) – date of the first data point in the data
- **data\_end** (*datetime.datetime*) – date of the last data point in the data
- **sim\_begin** (*datetime.datetime*) – date at which the simulation begins
- **sim\_end** (*datetime.datetime*) – date at which the simulation ends (should match fcast\_end)

- **fcst\_begin** (*datetime.datetime*) – date at which the forecast starts (should be one day after data\_end)
- **fcst\_end** (*datetime.datetime*) – data at which the forecast ends
- **data\_len** (*int*) – total number of days in the data
- **sim\_len** (*int*) – total number of days in the simulation
- **fcst\_len** (*int*) – total number of days in the forecast
- **diff\_data\_sim** (*int*) – difference in days between the simulation begin and the data begin. The simulation starting time is usually earlier than the data begin.

### Example

```
with Cov19Model(**params) as model:
    # Define model here
```

covid19\_inference.model.modelcontext (model)

return the given model or try to find it in the context if there was none supplied.

```
covid19_inference.model.student_t_likelihood(new_cases_inferred,
                                             pr_beta_sigma_obs=30,
                                             nu=4,
                                             offset_sigma=1,
                                             model=None,
                                             data_obs=None,
                                             name_student_t='_new_cases_studentT',
                                             name_sigma_obs='sigma_obs')
```

Set the likelihood to apply to the model observations (*model.new\_cases\_obs*) We assume a `StudentT` distribution because it is robust against outliers [Lange1989]. The likelihood follows:

$$P(\text{data\_obs}) \sim \text{StudentT}(\mu = \text{new\_cases\_inferred}, \sigma = \sigma_r, \text{nu} = \text{nu})$$

$$\sigma = \sigma_r \sqrt{\text{new\_cases\_inferred} + \text{offset\_sigma}}$$

The parameter  $\sigma_r$  follows a `HalfCauchy` prior distribution with parameter beta set by *pr\_beta\_sigma\_obs*. If the input is 2 dimensional, the parameter  $\sigma_r$  is different for every region.

#### Parameters

- **new\_cases\_inferred** (`TensorVariable`) – One or two dimensional array. If 2 dimensional, the first dimension is time and the second are the regions/countries
- **pr\_beta\_sigma\_obs** (*float*) – The beta of the `HalfCauchy` prior distribution of  $\sigma_r$ .
- **nu** (*float*) – How flat the tail of the distribution is. Larger nu should make the model more robust to outliers. Defaults to 4 [Lange1989].
- **offset\_sigma** (*float*) – An offset added to the sigma, to make the inference procedure robust. Otherwise numbers of *new\_cases\_inferred* would lead to very small errors and diverging likelihoods. Defaults to 1.
- **model** – The model on which we want to add the distribution
- **data\_obs** (*array*) – The data that is observed. By default it is *model.new\_cases\_obs*
- **name\_student\_t** – The name under which the studentT distribution is saved in the trace.
- **name\_sigma\_obs** – The name under which the distribution of the observable error is saved in the trace

**Returns** *None*

## References

`covid19_inference.model.SIR` (*lambda\_t\_log*, *mu*, *pr\_I\_begin=100*, *model=None*, *return\_all=False*, *save\_all=False*)

Implements the susceptible-infected-recovered model.

$$I_{new}(t) = \lambda_t I(t-1) \frac{S(t-1)}{N}$$

$$S(t) = S(t-1) - I_{new}(t)$$

$$I(t) = I(t-1) + I_{new}(t) - \mu I(t)$$

The prior distribution of the recovery rate  $\mu$  is set to *LogNormal*( $\log(\text{pr\_median\_mu})$ ),  $\text{pr\_sigma\_mu}$ ). And the prior distribution of  $I(0)$  to *HalfCauchy*( $\text{pr\_beta\_I\_begin}$ )

## Parameters

- **lambda\_t\_log** (*TensorVariable*) – time series of the logarithm of the spreading rate, 1 or 2-dimensional. If 2-dimensional the first dimension is time.
- **mu** (*TensorVariable*) – the recovery rate  $\mu$ , typically a random variable. Can be 0 or 1-dimensional. If 1-dimensional, the dimension are the different regions.
- **pr\_I\_begin** (float or array\_like or *TensorVariable*) – Prior beta of the Half-Cauchy distribution of  $I(0)$ .
- **pr\_median\_mu** (float or array\_like) – Prior for the median of the lognormal distribution of the recovery rate  $\mu$ .
- **pr\_sigma\_mu** (float or array\_like) – Prior for the sigma of the lognormal distribution of recovery rate  $\mu$ .
- **model** (*Cov19Model*) – if none, it is retrieved from the context
- **return\_all** (bool) – if True, returns *new\_I\_t*, *I\_t*, *S\_t* otherwise returns only *new\_I\_t*
- **save\_all** (bool) – if True, saves *new\_I\_t*, *I\_t*, *S\_t* in the trace, otherwise it saves only *new\_I\_t*

## Returns

- **new\_I\_t** (*TensorVariable*) – time series of the number daily newly infected persons.
- **I\_t** (*TensorVariable*) – time series of the infected (if *return\_all* set to True)
- **S\_t** (*TensorVariable*) – time series of the susceptible (if *return\_all* set to True)

`covid19_inference.model.SEIR` (*lambda\_t\_log*, *pr\_beta\_I\_begin=100*, *pr\_beta\_new\_E\_begin=50*, *pr\_median\_mu=0.125*, *pr\_mean\_median\_incubation=4*, *pr\_sigma\_median\_incubation=1*, *sigma\_incubation=0.4*, *pr\_sigma\_mu=0.2*, *model=None*, *return\_all=False*, *save\_all=False*, *name\_median\_incubation='median\_incubation'*)

Implements a model similar to the susceptible-exposed-infected-recovered model. Instead of a exponential decaying incubation period, the length of the period is lognormal distributed. The complete equation is:

$$\begin{aligned}
 E_{\text{new}}(t) &= \lambda_t I(t-1) \frac{S(t)}{N} \\
 S(t) &= S(t-1) - E_{\text{new}}(t) \\
 I_{\text{new}}(t) &= \sum_{k=1}^{10} \beta(k) E_{\text{new}}(t-k) \\
 I(t) &= I(t-1) + I_{\text{new}}(t) - \mu I(t) \\
 \beta(k) &= P(k) \sim \text{LogNormal}(\log(d_{\text{incubation}})), \text{sigma\_incubation})
 \end{aligned}$$

The recovery rate  $\mu$  and the incubation period is the same for all regions and follow respectively:

$$\begin{aligned}
 P(\mu) &\sim \text{LogNormal}(\log(\text{pr\_median\_mu}), \text{pr\_sigma\_mu}) \\
 P(d_{\text{incubation}}) &\sim \text{Normal}(\text{pr\_mean\_median\_incubation}, \text{pr\_sigma\_median\_incubation})
 \end{aligned}$$

The initial number of infected and newly exposed differ for each region and follow prior `HalfCauchy` distributions:

$$\begin{aligned}
 E(t) &\sim \text{HalfCauchy}(\text{pr\_beta\_E\_begin}) \quad \text{for } t \in \{-9, -8, \dots, 0\} \\
 I(0) &\sim \text{HalfCauchy}(\text{pr\_beta\_I\_begin}).
 \end{aligned}$$

### Parameters

- **lambda\_t\_log** (`TensorVariable`) – time series of the logarithm of the spreading rate, 1 or 2-dimensional. If 2-dimensional, the first dimension is time.
- **pr\_beta\_I\_begin** (`float` or `array_like`) – Prior beta of the `HalfCauchy` distribution of  $I(0)$ .
- **pr\_beta\_new\_E\_begin** (`float` or `array_like`) – Prior beta of the `HalfCauchy` distribution of  $E(0)$ .
- **pr\_median\_mu** (`float` or `array_like`) – Prior for the median of the `Lognormal` distribution of the recovery rate  $\mu$ .
- **pr\_mean\_median\_incubation** – Prior mean of the `Normal` distribution of the median incubation delay  $d_{\text{incubation}}$ . Defaults to 4 days [Nishiura2020], which is the median serial interval (the important measure here is not exactly the incubation period, but the delay until a person becomes infectious which seems to be about 1 day earlier as showing symptoms).
- **pr\_sigma\_median\_incubation** – Prior sigma of the `Normal` distribution of the median incubation delay  $d_{\text{incubation}}$ . Default is 1 day.
- **sigma\_incubation** – Scale parameter of the `Lognormal` distribution of the incubation time/ delay until infectiousness. The default is set to 0.4, which is about the scale found in [Nishiura2020], [Lauer2020].
- **pr\_sigma\_mu** (`float` or `array_like`) – Prior for the sigma of the lognormal distribution of recovery rate  $\mu$ .
- **model** (`Cov19Model`) – if none, it is retrieved from the context
- **return\_all** (`bool`) – if True, returns `new_I_t`, `new_E_t`, `I_t`, `S_t` otherwise returns only `new_I_t`
- **save\_all** (`bool`) – if True, saves `new_I_t`, `new_E_t`, `I_t`, `S_t` in the trace, otherwise it saves only `new_I_t`

- **name\_median\_incubation** (*str*) – The name under which the median incubation time is saved in the trace

### Returns

- **new\_I\_t** (*TensorVariable*) – time series of the number daily newly infected persons.
- **new\_E\_t** (*TensorVariable*) – time series of the number daily newly exposed persons. (if return\_all set to True)
- **I\_t** (*TensorVariable*) – time series of the infected (if return\_all set to True)
- **S\_t** (*TensorVariable*) – time series of the susceptible (if return\_all set to True)

### References

```
covid19_inference.model.delay_cases (new_I_t, pr_median_delay=10,
                                       pr_sigma_median_delay=0.2,
                                       pr_median_scale_delay=0.3,
                                       pr_sigma_scale_delay=None, model=None,
                                       save_in_trace=True, name_delay='delay',
                                       name_delayed_cases='new_cases_raw',
                                       len_input_arr=None, len_output_arr=None,
                                       diff_input_output=None)
```

Convolves the input by a lognormal distribution, in order to model a delay:

$$y_{\text{delayed}}(t) = \sum_{\tau=0}^T y_{\text{input}}(\tau) \text{LogNormal}[\log(\text{delay}), \text{pr\_median\_scale\_delay}](t - \tau)$$

$$\log(\text{delay}) = \text{Normal}(\log(\text{pr\_sigma\_delay}), \text{pr\_sigma\_delay})$$

For clarification: the *LogNormal* distribution is a function evaluated at  $t - \tau$ .

If the model is 2-dimensional, the  $\log(\text{delay})$  is hierarchically modelled with the *hierarchical\_normal()* function using the default parameters except that the prior  $\sigma$  of  $\text{delay}_{L2}$  is HalfNormal distributed (*error\_cauchy=False*).

### Parameters

- **new\_I\_t** (*TensorVariable*) – The input, typically the number newly infected cases  $I_{\text{new}}(t)$  of from the output of *SIR()* or *SEIR()*.
- **pr\_median\_delay** (*float*) – The mean of the normal distribution which models the prior median of the LogNormal delay kernel.
- **pr\_sigma\_median\_delay** (*float*) – The standart devaiation of normal distribution which models the prior median of the LogNormal delay kernel.
- **pr\_median\_scale\_delay** (*float*) – The scale (width) of the LogNormal delay kernel.
- **pr\_sigma\_scale\_delay** (*float*) – If it is not None, the scale is of the delay is kernel follows a prior LogNormal distribution, with median *pr\_median\_scale\_delay* and scale *pr\_sigma\_scale\_delay*.
- **model** (*Cov19Model*) – if none, it is retrieved from the context
- **save\_in\_trace** (*bool*) – whether to save  $y_{\text{delayed}}$  in the trace
- **name\_delay** (*str*) – The name under which the delay is saved in the trace, suffixes and prefixes are added depending on which variable is saved.

- **name\_delayed\_cases** (*str*) – The name under which the delay is saved in the trace, suffixes and prefixes are added depending on which variable is saved.
- **len\_input\_arr** – Length of `new_I_t`. By default equal to `model.sim_len`. Necessary because the shape of theano tensors are not defined at when the graph is built.
- **len\_output\_arr** (*int*) – Length of the array returned. By default it set to the length of the `cases_obs` saved in the model plus the number of days of the forecast.
- **diff\_input\_output** (*int*) – Number of days the returned array begins later then the input. Should be significantly larger than the median delay. By default it is set to the `model.diff_data_sim`.

**Returns** `new_cases_inferred` (*TensorVariable*) – The delayed input  $y_{\text{delayed}}(t)$ , typically the daily number new cases that one expects to measure.

```
covid19_inference.model.week_modulation(new_cases_raw, week_modulation_type='abs_sine',
                                         pr_mean_weekend_factor=0.3,
                                         pr_sigma_weekend_factor=0.5,
                                         week_end_days=(6, 7), model=None,
                                         save_in_trace=True)
```

Adds a weekly modulation of the number of new cases:

$$\text{new\_cases} = \text{new\_cases\_raw} \cdot (1 - f(t)), \quad \text{with}$$

$$f(t) = f_w \cdot \left(1 - \left|\sin\left(\frac{\pi}{7}t - \frac{1}{2}\Phi_w\right)\right|\right),$$

if `week_modulation_type` is "abs\_sine" (the default). If `week_modulation_type` is "step", the new cases are simply multiplied by the weekend factor on the days set by `week_end_days`

The weekend factor  $f_w$  follows a Lognormal distribution with median `pr_mean_weekend_factor` and sigma `pr_sigma_weekend_factor`. It is hierarchically constructed if the input is two-dimensional by the function `hierarchical_normal()` with default arguments.

The offset from Sunday  $\Phi_w$  follows a flat `VonMises` distribution and is the same for all regions.

#### Parameters

- **new\_cases\_raw** (*TensorVariable*) – The input array, can be one- or two-dimensional
- **week\_modulation\_type** (*str*) – The type of modulation, accepts "step" or "abs\_sine" (the default).
- **pr\_mean\_weekend\_factor** (*float*) – Sets the prior mean of the factor  $f_w$  by which weekends are counted.
- **pr\_sigma\_weekend\_factor** (*float*) – Sets the prior sigma of the factor  $f_w$  by which weekends are counted.
- **week\_end\_days** (*tuple of ints*) – The days counted as weekend if `week_modulation_type` is "step"
- **model** (*Cov19Model*) – if none, it is retrieved from the context
- **save\_in\_trace** (*bool*) – If True (default) the new\_cases are saved in the trace.

**Returns** `new_cases` (*TensorVariable*)

```
covid19_inference.model.make_change_point_RVs(change_points_list,
                                                pr_median_lambda_0,
                                                pr_sigma_lambda_0=1, model=None)
```

#### Parameters



- **priors\_dict** –
- **change\_points\_list** –
- **model** –

```
covid19_inference.model.lambda_t_with_sigmoids(change_points_list,
                                                pr_median_lambda_0,
                                                pr_sigma_lambda_0=0.5,
                                                model=None)
```

#### Parameters

- **change\_points\_list** –
- **pr\_median\_lambda\_0** –
- **pr\_sigma\_lambda\_0** –
- **model** (*Cov19Model*) – if none, it is retrieved from the context

```
covid19_inference.model.hierarchical_normal(name, name_sigma, pr_mean, pr_sigma,
                                             len_L2, w=1.0, error_fact=2.0, error_cauchy=True)
```

Implements an hierarchical normal model:

$$\begin{aligned}x_{L1} &= \text{Normal}(\text{pr\_mean}, \text{pr\_sigma}) \\ y_{i,L2} &= \text{Normal}(x_{L1}, \sigma_{L2}) \\ \sigma_{L2} &= \text{HalfCauchy}(\text{error\_fact} \cdot \text{pr\_sigma})\end{aligned}$$

It is however implemented in a non-centered way, that the second line is changed to:

$$y_{i,L2} = x_{L1} + \text{Normal}(0, 1) \cdot \sigma_{L2}$$

See for example <https://arxiv.org/pdf/1312.0906.pdf>

#### Parameters

- **name** (*str*) – Name under which  $x_{L1}$  and  $y_{L2}$  saved in the trace. `'_L1'` and `'_L2'` is appended
- **name\_sigma** (*str*) – Name under which  $\sigma_{L2}$  saved in the trace. `'_L2'` is appended.
- **pr\_mean** (*float*) – Prior mean of  $x_{L1}$
- **pr\_sigma** (*float*) – Prior sigma for  $x_{L1}$  and (multiplied by `error_fact`) for  $\sigma_{L2}$
- **len\_L2** (*int*) – length of  $y_{L2}$
- **error\_fact** (*float*) – Factor by which `pr_sigma` is multiplied as prior for `sigma_text{L2}`
- **error\_cauchy** (*bool*) – if False, a *HalfNormal* distribution is used for  $\sigma_{L2}$  instead of *HalfCauchy*

#### Returns

- **y** (*TensorVariable*) – the random variable  $y_{L2}$
- **x** (*TensorVariable*) – the random variable  $x_{L1}$

```
covid19_inference.model.make_prior_I(lambda_t_log, mu, pr_median_delay,  
                                     pr_sigma_I_begin=2, n_data_points_used=5,  
                                     model=None)
```

Builds the prior for I begin by solving the SIR differential from the first data backwards. This decorrelates the I\_begin from the lambda\_t at the beginning, allowing a more efficient sampling. The example\_one\_bundesland runs about 30% faster with this prior, instead of a HalfCauchy.

**Parameters**

- **lambda\_t\_log** (*TensorVariable*) –
- **mu** (*TensorVariable*) –
- **pr\_median\_delay** (*float*) –
- **pr\_sigma\_I\_begin** (*float*) –
- **n\_data\_points\_used** (*int*) –
- **model** (*Cov19Model*) – if none, it is retrieved from the context

**Returns** **I\_begin** (*TensorVariable*)

## DATA RETRIEVAL

### Table of Contents

- *Data Retrieval*
  - *Utility*
  - *Johns Hops University*
  - *Robert Koch Institute*
  - *Robert Koch Institute situation reports*
  - *Google*
  - *Our World in Data*
  - *Base Retrieval Class*

## 5.1 Utility

`covid19_inference.data_retrieval.retrieval.set_data_dir` (*fname=None*, *permissions=None*)

Set the global variable `_data_dir`. New downloaded data is placed there. If no argument provided we try the default tmp directory. If permissions are not provided, uses defaults if *fname* is in user folder. If not in user folder, tries to set 777.

`covid19_inference.data_retrieval.retrieval.backup_instances` (*trace=None*,  
*model=None*,  
*fname='latest\_'*)

helper to save or load trace and model instances. loads from *fname* if provided traces and model variables are None, else saves them there.

## 5.2 Johns Hops University

**class** covid19\_inference.data\_retrieval.JHU(*auto\_download=False*)

This class can be used to retrieve and filter the dataset from the online repository of the coronavirus visual dashboard operated by the [Johns Hopkins University](#).

### Features

- download all files from the online repository of the coronavirus visual dashboard operated by the Johns Hopkins University.
- filter by deaths, confirmed cases and recovered cases
- filter by country and state
- filter by date

### Example

```
jhu = cov19.data_retrieval.JHU()
jhu.download_all_available_data()

#Access the data by
jhu.data
#or
jhu.get_new("confirmed", "Italy")
jhu.get_total(filter)
```

**\_\_init\_\_**(*auto\_download=False*)

On init of this class the base Retrieval Class **\_\_init\_\_** is called, with jhu specific arguments.

**Parameters** *auto\_download* (*bool*, *optional*) – Whether or not to automatically call the `download_all_available_data()` method. One should explicitly call this method for more configuration options (default: `false`)

**download\_all\_available\_data**(*force\_local=False*, *force\_download=False*)

Attempts to download from the main urls (`self.url_csv`) which was set on initialization of this class. If this fails it downloads from the fallbacks. It can also be specified to use the local files or to force the download. The download methods get inherited from the base retrieval class.

### Parameters

- **force\_local** (*bool*, *optional*) – If True forces to load the local files.
- **force\_download** (*bool*, *optional*) – If True forces the download of new files

**get\_total\_confirmed\_deaths\_recovered**(*country: str = None*, *state: str = None*, *begin\_date: datetime.datetime = None*, *end\_date: datetime.datetime = None*)

Retrieves all confirmed, deaths and recovered cases from the Johns Hopkins University dataset as a DataFrame with datetime index. Can be filtered by country and state, if only a country is given all available states get summed up.

### Parameters

- **country** (*str*, *optional*) – name of the country (the “Country/Region” column), can be None if the whole summed up data is wanted (why would you do this?)
- **state** (*str*, *optional*) – name of the state (the “Province/State” column), can be None if country is set or the whole summed up data is wanted

- **begin\_date** (*datetime.datetime, optional*) – initial date for the returned data, if no value is given the first date in the dataset is used
- **end\_date** (*datetime.datetime, optional*) – last date for the returned data, if no value is given the most recent date in the dataset is used

**Returns** *pandas.DataFrame*

**get\_new** (*value='confirmed', country: str = None, state: str = None, data\_begin: datetime.datetime = None, data\_end: datetime.datetime = None*)

Retrieves all new cases from the Johns Hopkins University dataset as a DataFrame with datetime index. Can be filtered by value, country and state, if only a country is given all available states get summed up.

#### Parameters

- **value** (*str*) – Which data to return, possible values are - “confirmed”, - “recovered”, - “deaths” (default: “confirmed”)
- **country** (*str, optional*) – name of the country (the “Country/Region” column), can be None
- **state** (*str, optional*) – name of the state (the “Province/State” column), can be None
- **begin\_date** (*datetime.datetime, optional*) – initial date for the returned data, if no value is given the first date in the dataset is used
- **end\_date** (*datetime.datetime, optional*) – last date for the returned data, if no value is given the most recent date in the dataset is used

**Returns** *pandas.DataFrame* – table with new cases and the date as index

**get\_total** (*value='confirmed', country: str = None, state: str = None, data\_begin: datetime.datetime = None, data\_end: datetime.datetime = None*)

Retrieves all total/cumulative cases from the Johns Hopkins University dataset as a DataFrame with date-time index. Can be filtered by value, country and state, if only a country is given all available states get summed up.

#### Parameters

- **value** (*str*) – Which data to return, possible values are - “confirmed”, - “recovered”, - “deaths” (default: “confirmed”)
- **country** (*str, optional*) – name of the country (the “Country/Region” column), can be None
- **state** (*str, optional*) – name of the state (the “Province/State” column), can be None
- **begin\_date** (*datetime.datetime, optional*) – initial date for the returned data, if no value is given the first date in the dataset is used
- **end\_date** (*datetime.datetime, optional*) – last date for the returned data, if no value is given the most recent date in the dataset is used

**Returns** *pandas.DataFrame* – table with total/cumulative cases and the date as index

**filter\_date** (*df, begin\_date: datetime.datetime = None, end\_date: datetime.datetime = None*)

Returns give dataframe between begin and end date. Dataframe has to have a datetime index.

#### Parameters

- **begin\_date** (*datetime.datetime, optional*) – First day that should be filtered

- **end\_date** (*datetime.datetime, optional*) – Last day that should be filtered

Returns *pandas.DataFrame*

**get\_possible\_countries\_states()**

Can be used to get a list with all possible states and countries.

Returns *pandas.DataFrame* in the format

## 5.3 Robert Koch Institute

**class** covid19\_inference.data\_retrieval.**RKI** (*auto\_download=False*)

This class can be used to retrieve and filter the dataset from the Robert Koch Institute [Robert Koch Institute](#). The data gets retrieved from the [arcgis](#) dashboard.

### Features

- download the full dataset
- filter by date
- filter by bundesland
- filter by recovered, deaths and confirmed cases

### Example

```
rki = covid19_inference.data_retrieval.RKI()
rki.download_all_available_data()

#Access the data by
rki.data
#or
rki.get_new("confirmed", "Sachsen")
rki.get_total(filter)
```

**\_\_init\_\_** (*auto\_download=False*)

On init of this class the base Retrieval Class **\_\_init\_\_** is called, with rki specific arguments.

**Parameters** **auto\_download** (*bool, optional*) – Whether or not to automatically call the `download_all_available_data()` method. One should explicitly call this method for more configuration options (default: false)

**download\_all\_available\_data** (*force\_local=False, force\_download=False*)

Attempts to download from the main url (self.url\_csv) which was given on initialization. If this fails download from the fallbacks. It can also be specified to use the local files or to force the download. The download methods get inherited from the base retrieval class.

### Parameters

- **force\_local** (*bool, optional*) – If True forces to load the local files.
- **force\_download** (*bool, optional*) – If True forces the download of new files

**get\_total** (*value='confirmed', bundesland: str = None, landkreis: str = None, data\_begin: datetime.datetime = None, data\_end: datetime.datetime = None, date\_type: str = 'date'*)

Gets all total confirmed cases for a region as dataframe with date index. Can be filtered with multiple arguments.

### Parameters

- **value** (*str*) – Which data to return, possible values are - “confirmed”, - “recovered”, - “deaths” (default: “confirmed”)
- **bundesland** (*str*, *optional*) – if no value is provided it will use the full summed up dataset for Germany
- **landkreis** (*str*, *optional*) – if no value is provided it will use the full summed up dataset for the region (bundesland)
- **data\_begin** (*datetime.datetime*, *optional*) – initial date, if no value is provided it will use the first possible date
- **data\_end** (*datetime.datetime*, *optional*) – last date, if no value is provided it will use the most recent possible date
- **date\_type** (*str*, *optional*) – type of date to use: reported date ‘date’ (Meldedatum in the original dataset), or symptom date ‘date\_ref’ (Refdatum in the original dataset)

**Returns** *pandas.DataFrame*

**get\_new** (*value*='confirmed', *bundesland*: *str* = None, *landkreis*: *str* = None, *data\_begin*: *datetime.datetime* = None, *data\_end*: *datetime.datetime* = None, *date\_type*: *str* = 'date')

Retrieves all new cases from the Robert Koch Institute dataset as a DataFrame with datetime index. Can be filtered by value, bundesland and landkreis, if only a country is given all available states get summed up.

#### Parameters

- **value** (*str*) – Which data to return, possible values are - “confirmed”, - “recovered”, - “deaths” (default: “confirmed”)
- **bundesland** (*str*, *optional*) – if no value is provided it will use the full summed up dataset for Germany
- **landkreis** (*str*, *optional*) – if no value is provided it will use the full summed up dataset for the region (bundesland)
- **data\_begin** (*datetime.datetime*, *optional*) – initial date for the returned data, if no value is given the first date in the dataset is used, if none is given could yield errors
- **data\_end** (*datetime.datetime*, *optional*) – last date for the returned data, if no value is given the most recent date in the dataset is used

**Returns** *pandas.DataFrame* – table with daily new confirmed and the date as index

**filter** (*data\_begin*: *datetime.datetime* = None, *data\_end*: *datetime.datetime* = None, *variable*='confirmed', *date\_type*='date', *level*=None, *value*=None)

Filters the obtained dataset for a given time period and returns an array ONLY containing only the desired variable.

#### Parameters

- **data\_begin** (*datetime.datetime*, *optional*) – initial date, if no value is provided it will use the first possible date
- **data\_end** (*datetime.datetime*, *optional*) – last date, if no value is provided it will use the most recent possible date
- **variable** (*str*, *optional*) – type of variable to return possible types are: “confirmed”: cases (default) “AnzahlTodesfall”: deaths “AnzahlGenesen”: recovered
- **date\_type** (*str*, *optional*) – type of date to use: reported date ‘date’ (Meldedatum in the original dataset), or symptom date ‘date\_ref’ (Refdatum in the original dataset)

- **level** (*str*, *optional*) –

**possible strings are:** "None" : return data from all Germany (default) "Bundesland" : a state "Landkreis" : a region

- **value** (*None*, *optional*) – string of the state/region e.g. "Sachsen"

**Returns** *pd.DataFrame* – array with ONLY the requested variable, in the requested range. (one dimensional)

**filter\_all\_bundesland** (*begin\_date: datetime.datetime = None, end\_date: datetime.datetime = None, variable='confirmed', date\_type='date'*)

Filters the full RKI dataset

#### Parameters

- **df** (*DataFrame*) – RKI dataframe, from `get_rki()`
- **begin\_date** (*datetime.datetime*) – initial date to return
- **end\_date** (*datetime.datetime*) – last date to return
- **variable** (*str*, *optional*) – type of variable to return: cases ("AnzahlFall"), deaths ("AnzahlTodesfall"), recovered ("AnzahlGenesen")
- **date\_type** (*str*, *optional*) – type of date to use: reported date 'date' (Meldedatum in the original dataset), or symptom date 'date\_ref' (Refdatum in the original dataset)

**Returns** *pd.DataFrame* – DataFrame with datetime dates as index, and all German regions (bundesländer) as columns

## 5.4 Robert Koch Institute situation reports

**class** `covid19_inference.data_retrieval.RKIsituationreports` (*auto\_download=False*)

As mentioned by Matthias Linden, the daily situation reports have more available data. This class retrieves this additional data from Matthias website and parses it into the format we use i.e. a datetime index.

Interesting new data is for example ICU cases, deaths and recorded symptoms. For now one can look at the data by running

#### Example

```
rki_si_re = covid19.data_retrieval.RKIsituationreports(True)
print(rki_si_re.data)
```

---

**Todo:** Filter functions for ICU, Symptoms and maybe even daily new cases for the respective categories.

---

**\_\_init\_\_** (*auto\_download=False*)

On init of this class the base Retrieval Class `__init__` is called, with rki situation reports specific arguments.

**Parameters** **auto\_download** (*bool*, *optional*) – Whether or not to automatically call the `download_all_available_data()` method. One should explicitly call this method for more configuration options (default: false)

**download\_all\_available\_data** (*force\_local=False, force\_download=False*)

Attempts to download from the main url (self.url\_csv) which was given on initialization. If this fails



download from the fallbacks. It can also be specified to use the local files or to force the download. The download methods get inherited from the base retrieval class.

#### Parameters

- **force\_local** (*bool*, *optional*) – If True forces to load the local files.
- **force\_download** (*bool*, *optional*) – If True forces the download of new files

## 5.5 Google

**class** covid19\_inference.data\_retrieval.GOOGLE (*auto\_download=False*)

This class can be used to retrieve the mobility dataset from [Google](#).

#### Example

```
gl = covid19.data_retrieval.GOOGLE()
gl.download_all_available_data()

#Access the data by
gl.data
#or
gl.get_changes(filter)
```

**\_\_init\_\_** (*auto\_download=False*)

On init of this class the base Retrieval Class **\_\_init\_\_** is called, with google specific arguments.

**Parameters** **auto\_download** (*bool*, *optional*) – Whether or not to automatically call the `download_all_available_data()` method. One should explicitly call this method for more configuration options (default: false)

**download\_all\_available\_data** (*force\_local=False*, *force\_download=False*)

Attempts to download from the main url (self.url\_csv) which was given on initialization. If this fails download from the fallbacks. It can also be specified to use the local files or to force the download. The download methods get inherited from the base retrieval class.

#### Parameters

- **force\_local** (*bool*, *optional*) – If True forces to load the local files.
- **force\_download** (*bool*, *optional*) – If True forces the download of new files

**get\_changes** (*country: str*, *state: str = None*, *region: str = None*, *data\_begin: datetime.datetime = None*, *data\_end: datetime.datetime = None*)

Returns a dataframe with the relative changes in mobility to a baseline, provided by google. They are separated into “retail and recreation”, “grocery and pharmacy”, “parks”, “transit”, “workplaces” and “residential”. Filterable for country, state and region and date.

#### Parameters

- **country** (*str*) – Selected country for the mobility data.
- **state** (*str*, *optional*) – State for the selected data if no value is selected the whole country is chosen
- **region** (*str*, *optional*) – Region for the selected data if no value is selected the whole region/country is chosen
- **data\_end** (*data\_begin*,) – Filter for the desired time period

**Returns** *pandas.DataFrame*

**get\_possible\_counties\_states\_regions()**

Can be used to obtain all different possible countries with there corresponding possible states and regions.

**Returns** *pandas.DataFrame*

## 5.6 Our World in Data

**class** covid19\_inference.data\_retrieval.OWD(*auto\_download=False*)

This class can be used to retrieve the testings dataset from [Our World in Data](#).

### Example

```
owd = covid19_inference.data_retrieval.OWD()
owd.download_all_available_data()
```

**\_\_init\_\_**(*auto\_download=False*)

On init of this class the base Retrieval Class **\_\_init\_\_** is called, with google specific arguments.

**Parameters** **auto\_download** (*bool, optional*) – Whether or not to automatically call the **download\_all\_available\_data()** method. One should explicitly call this method for more configuration options (default: *false*)

**download\_all\_available\_data**(*force\_local=False, force\_download=False*)

Attempts to download from the main url (self.url\_csv) which was given on initialization. If this fails download from the fallbacks. It can also be specified to use the local files or to force the download. The download methods get inhereted from the base retrieval class.

#### Parameters

- **force\_local** (*bool, optional*) – If True forces to load the local files.
- **force\_download** (*bool, optional*) – If True forces the download of new files

**get\_possible\_countries()**

Can be used to obtain all different possible countries in the dataset.

**Returns** *pandas.DataFrame*

**get\_total**(*value='tests', country=None, data\_begin=None, data\_end=None*)

Retrieves all new cases from the Our World in Data dataset as a DataFrame with datetime index. Can be filtered by value, country and state, if only a country is given all available states get summed up.

#### Parameters

- **value** (*str*) – Which data to return, possible values are - “confirmed”, - “tests”, - “deaths” (default: “confirmed”)
- **country** (*str*) – name of the country
- **begin\_date** (*datetime.datetime, optional*) – initial date for the returned data, if no value is given the first date in the dataset is used
- **end\_date** (*datetime.datetime, optional*) – last date for the returned data, if no value is given the most recent date in the dataset is used

**Returns** *pandas.DataFrame* – table with new cases and the date as index

**get\_new** (*value='tests', country=None, data\_begin=None, data\_end=None*)

Retrieves all new cases from the Our World in Data dataset as a DataFrame with datetime index. casesn be filtered by value, country and state, if only a country is given all available states get summed up.

#### Parameters

- **value** (*str*) – Which data to return, possible values are - “confirmed”, - “tests”, - “deaths” (default: “confirmed”)
- **country** (*str*) – name of the country
- **begin\_date** (*datetime.datetime, optional*) – initial date for the returned data, if no value is given the first date in the dataset is used
- **end\_date** (*datetime.datetime, optional*) – last date for the returned data, if no value is given the most recent date in the dataset is used

**Returns** *pandas.DataFrame* – table with new cases and the date as index

## 5.7 Base Retrieval Class

```
class covid19_inference.data_retrieval.retrieval.Retrieval (name, url_csv,  
 fallbacks, up-  
 date_interval=None,  
 **kwargs)
```

Each source class should inherit this base retrieval class, it streamlines alot of base functions. It manages downloads, multiple fallbacks and local backups via timestamp. At init of the parent class the Retrieval init should be called with the following arguments, these get saved as attributes.

An example for the usage can be seen in the \_Google, \_RKI and \_JHU source files.

```
__init__ (name, url_csv, fallbacks, update_interval=None, **kwargs)
```

#### Parameters

- **name** (*str*) – A name for the Parent class, mainly used for the local file backup.
- **url\_csv** (*str*) – The url to the main dataset as csv, if an empty string if supplied the fallback routines get used.
- **fallbacks** (*array*) – Fallbacks can be filepaths to local or online sources or even methods defined in the parent class.
- **update\_interval** (*datetime.timedelta*) – If the local file is older than the update\_interval it gets updated once the download all function is called.

```
_download_csv_from_source (filepath, **kwargs)
```

Uses pandas read csv to download the csv file. The possible kwargs can be seen in the pandas [documentation](#).

These kwargs can vary for the different parent classes and should be defined there!

**filepath** [str] Full path to the desired csv file

**Returns** *bool* – True if the retrieval was a success, False if it failed

```
_fallback_handler ()
```

Recursively iterate over all fallbacks and try to execute subroutines depending on the type of fallback.

```
_timestamp_local_old (force_local=False) → bool
```

1. Get timestamp if it exists
2. compare with the date today
3. update if data is older than set intervall -> can be parent dependant

**`_save_to_local()`**

Creates a local backup for the self.data pandas.DataFrame. And a timestamp for the source.

## PLOTTING

`covid19_inference.plotting.get_all_free_RVs_names(model)`

Returns the names of all free parameters of the model

**Parameters** `model` (*pm.Model instance*) –

**Returns** *list* – all variable names

`covid19_inference.plotting.get_prior_distribution(model, x, varname)`

Given a model and variable name, get the prior that was used for modeling.

**Parameters**

- **model** (*pm.Model instance*) –

- **x** (*list or array*) –

- **varname** (*string*) –

**Returns** *array* – the prior distribution evaluated at x

`covid19_inference.plotting.plot_hist(model, trace, ax, varname, colors=('tab:blue', 'tab:orange'), bins=50)`

Plots one histogram of the prior and posterior distribution of the variable varname.

**Parameters**

- **model** (*pm.Model instance*) –

- **trace** (*trace of the model*) –

- **ax** (*matplotlib.axes instance*) –

- **varname** (*string*) –

- **colors** (*list with 2 colornames*) –

- **bins** (*number or array*) – passed to np.hist

**Returns** *None*

`covid19_inference.plotting.plot_cases(trace, new_cases_obs, date_begin_sim, diff_data_sim, start_date_plot=None, end_date_plot=None, ylim=None, week_interval=None, colors=('tab:blue', 'tab:orange'), country='Germany')`

Plots the new cases, the fit, forecast and lambda\_t evolution

**Parameters**

- **trace** (*trace returned by model*) –

- **new\_cases\_obs** (*array*) –

- `date_begin_sim(datetime.datetime)` –
- `diff_data_sim(float)` – Difference in days between the begin of the simulation and the data
- `start_date_plot(datetime.datetime)` –
- `end_date_plot(datetime.datetime)` –
- `ylim(float)` – the maximal y value to be plotted
- `week_interval(int)` – the interval in weeks of the y ticks
- `colors(list with 2 colornames)` –

**Returns** *figure, axes*

## VARIABLES SAVED IN THE TRACE

The trace by default contains the following parameters in the SIR/SEIR hierarchical model. XXX denotes a number.

Name in trace	Dimensions	Created by function
lambda_XXX_L1	samples	lambda_t_with_sigmoids/make_change_point_RVs
lambda_XXX_L2	samples x re- gions	lambda_t_with_sigmoids/make_change_point_RVs
sigma_lambda_XXX_L1	samples	lambda_t_with_sigmoids/make_change_point_RVs
transient_day_XXX_L1	samples	lambda_t_with_sigmoids/make_change_point_RVs
transient_day_XXX_L2	samples x re- gions	lambda_t_with_sigmoids/make_change_point_RVs
sigma_transient_day_XXX_L2	samples	lambda_t_with_sigmoids/make_change_point_RVs
transient_len_XXX_L1	samples	lambda_t_with_sigmoids/make_change_point_RVs
transient_len_XXX_L2	samples x re- gions	lambda_t_with_sigmoids/make_change_point_RVs
sigma_transient_len_XXX_L2	samples	lambda_t_with_sigmoids/make_change_point_RVs
delay_L1	samples	delay_cases
delay_L2	samples x re- gions	delay_cases
sigma_delay_L2	samples	delay_cases
weekend_factor_L1	samples	week_modulation
weekend_factor_L2	samples x re- gions	week_modulation
sigma_weekend_factor_L2	samples	week_modulation
offset_modulation	samples	week_modulation
new_cases_raw	samples x time x regions	week_modulation
mu	samples	SIR/SEIR
I_begin	samples x re- gions	SIR/SEIR
new_cases	samples x time x regions	SIR/SEIR
sigma_obs	samples x re- gions	SIR/SEIR
new_E_begin	samples x 11 x regions	SEIR
median_incubation_start	samples	SEIR
median_incubation_start_L2	samples x re- gions	SEIR
sigma_median_incubation_start_L2	samples	SEIR

For the non-hierarchical model, variables with \_L2 suffixes are missing, and \_L1 suffixes are removed from the name.



## CONTRIBUTING

We always welcome contributions. Here we gather some guidelines to make the process as smooth as possible.

### 8.1 Beginning

To see where help is needed, go to the issues page on Github. If you want to begin on an issue, make a comment below and begin a draft pull request: <https://github.blog/2019-02-14-introducing-draft-pull-requests/> You can link the pull request on the right side of the commit to it.

When you have finished working on the issue, change it to a regular pull request. Check that there are no conflicts to the current master (<https://www.digitalocean.com/community/tutorials/how-to-rebase-and-update-a-pull-request>)

### 8.2 Code formatting

We use black <https://github.com/psf/black> as automatic code formatter. Please run your code through it before you open a pull request.

We do not check for formatting in the testing (travis) but recommend to set up black as a pre-commit hook.

```
conda install -c conda-forge pre-commit
pre-commit install
```

Try to stick to PEP 8. You can use [type annotations](#) if you want, but it is not necessary or encouraged.

### 8.3 Testing

We use travis and pytest. To check your changes locally:

```
python -m pytest --log-level=INFO --log-cli-level=INFO
```

It would be great if anything that is added to the code-base has an according test in the `tests` folder. We are not there yet, but it is on the todo. Be encouraged to add tests :)

## 8.4 Documentation

The documentation is built using Sphinx from the docstrings. To test it before submitting, navigate with a terminal to the docs/ directory. Install if necessary the packages listed in `piprequirements.txt` run `make html`. The documentation can then be accessed in `docs/_build/html/index.html`. As an example you can look at the documentation of `covid19_inference.model.SIR()`

## DEBUGGING

This is some pointer to help debugging models and sampling issues

### 9.1 General approach for nans/infs during sampling

The idea of this approach is to sample from the prior and then run the model. If the log likelihood is then  $-\infty$ , there is a problem, and the output of the theano functions is inspected.

Sample from prior:

```
from pymc3.util import (
    get_untransformed_name,
    is_transformed_name)

varnames = list(map(str, model.vars))

for name in varnames:
    if is_transformed_name(name):
        varnames.append(get_untransformed_name(name))

with model:
    points = pm.sample_prior_predictive(var_names = varnames)
    points_list = []
    for i in range(len(next(iter(points.values())))):
        point_dict = {}
        for name, val in points.items():
            point_dict[name] = val[i]
        points_list.append(point_dict)
```

points\_list is a list of the starting points for the model, sampled from the prior. Then to run the model and print the log-likelihood:

```
fn = model.fn(model.logpt)

for point in points_list[:]:
    print(fn(point))
```

To monitor the output and save it in a file (for use in ipython). Learned from: [http://deeplearning.net/software/theano/tutorial/debug\\_faq.html#how-do-i-step-through-a-compiled-function](http://deeplearning.net/software/theano/tutorial/debug_faq.html#how-do-i-step-through-a-compiled-function)

```
%%capture cap --no-stderr
def inspect_inputs(i, node, fn):
    print(i, node, "input(s) value(s):", [input[0] for input in fn.inputs],
```

(continues on next page)

(continued from previous page)

```

        end='')

def inspect_outputs(i, node, fn):
    print(" output(s) value(s):", [output[0] for output in fn.outputs])

fn_monitor = model.fn(model.logpt,
                       mode=theano.compile.MonitorMode(
                           pre_func=inspect_inputs,
                           post_func=inspect_outputs).excluding(
                               'local_elemwise_fusion', 'inplace'))

fn = model.fn(model.logpt)

for point in points_list[:]:
    if fn(point) < -1e10:
        print(fn_monitor(point))
        break

```

In a new cell:

```

with open('output.txt', 'w') as f:
    f.write(cap.stdout)

```

Then one can open output.txt in a text editor, and follow from where infs or nans come from by following the inputs and outputs up through the graph

## 9.2 Sampler: MCMC (Nuts)

### 9.2.1 Divergences

During sampling, a significant fraction of divergences are a sign that the sampler doesn't sample the whole posterior. In this case the model should be reparametrized. See this tutorial for a typical example: [https://docs.pymc.io/notebooks/Diagnosing\\_biased\\_Inference\\_with\\_Divergences.html](https://docs.pymc.io/notebooks/Diagnosing_biased_Inference_with_Divergences.html)

And these papers include some more details: <https://pdfs.semanticscholar.org/7b85/fb48a077c679c325433fbe13b87560e12886.pdf> <https://arxiv.org/pdf/1312.0906.pdf>

### 9.2.2 Bad initial energy

This typically occurs when some distribution in the model can't be evaluated at the starting point of chain. Run this to see which distribution throws nans or infs:

```

for RV in model.basic_RVs:
    print(RV.name, RV.logp(model.test_point))

```

However, this only evaluates the test\_point. When PyMC3 starts sampling, it adds some jitter around this test\_point, which then could lead to nans. Run this to add jitter and then evaluate the logp:

```

chains=4
for RV in model.basic_RVs:
    print(RV.name)

```

(continues on next page)

(continued from previous page)

```

for _ in range(chains):
    mean = {var: val.copy() for var, val in model.test_point.items()}
    for val in mean.values():
        val[...] += 2 * np.random.rand(*val.shape) - 1
    print(RV.logp(mean))

```

This code could potentially change in newer versions of PyMC3 (this is tested in 3.8). Read the source code, to know which random jitter PyMC3 currently adds at beginning.

### 9.2.3 Nans occur during sampling

Run the sampler with the debug mode of Theano.

```

from theano.compile.nanguardmode import NanGuardMode
mode = NanGuardMode(nan_is_error=True, inf_is_error=False, big_is_error=False,
                    optimizer='ol')
trace = pm.sample(mode=mode)

```

However this doesn't lead to helpful messages if nans occur during gradient evaluations.

## 9.3 Sampler: Variational Inference

There exist some ways to track parameters during sampling. An example:

```

with model:
    advi = pm.ADVI()
    print(advi.approx.group)

    print(advi.approx.mean.eval())
    print(advi.approx.std.eval())

    tracker = pm.callbacks.Tracker(
        mean=advi.approx.mean.eval, # callable that returns mean
        std=advi.approx.std.eval    # callable that returns std
    )

    approx = advi.fit(100000, callbacks=[tracker],
                     obj_optimizer=pm.adagrad_window(learning_rate=1e-3),
                     #total_grad_norm_constraint=10) #constrains maximal gradient,
    ↪could help

    print(approx.groups[0].bij.rmap(approx.params[0].eval()))

    plt.plot(tracker['mean'])
    plt.plot(tracker['std'])

```

For the tracker, the order of the parameters is saved in:

```
approx.ordering.by_name
```

and the indices encoded there in the slc field. To plot the mean value of a given parameter name, run:

```
plt.plot(np.array(tracker['mean']))[:, approx.ordering.by_name['name'].slc]
```

The debug mode is set with the following parameter:

```
from theano.compile.nanguardmode import NanGuardMode
mode = NanGuardMode(nan_is_error=True, inf_is_error=False, big_is_error=False,
                    optimizer='ol')
approx = advi.fit(100000, callbacks=[tracker],
                 fn_kwargs={'mode':mode})
```

## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`





## BIBLIOGRAPHY

- [Lange1989] Lange, K., Roderick J. A. Little, & Jeremy M. G. Taylor. (1989). Robust Statistical Modeling Using the  $t$  Distribution. *Journal of the American Statistical Association*, 84(408), 881-896. doi:10.2307/2290063
- [Nishiura2020] Nishiura, H.; Linton, N. M.; Akhmetzhanov, A. R. Serial Interval of Novel Coronavirus (COVID-19) Infections. *Int. J. Infect. Dis.* 2020, 93, 284–286. <https://doi.org/10.1016/j.ijid.2020.02.060>.
- [Lauer2020] Lauer, S. A.; Grantz, K. H.; Bi, Q.; Jones, F. K.; Zheng, Q.; Meredith, H. R.; Azman, A. S.; Reich, N. G.; Lessler, J. The Incubation Period of Coronavirus Disease 2019 (COVID-19) From Publicly Reported Confirmed Cases: Estimation and Application. *Ann Intern Med* 2020. <https://doi.org/10.7326/M20-0504>.



## PYTHON MODULE INDEX

### C

covid19\_inference, ??  
covid19\_inference.data\_retrieval.retrieval,  
15  
covid19\_inference.model, 7  
covid19\_inference.plotting, 25



# INDEX

## Symbols

`__init__()` (*covid19\_inference.data\_retrieval.GOOGLE* method), 21  
`__init__()` (*covid19\_inference.data\_retrieval.JHU* method), 16  
`__init__()` (*covid19\_inference.data\_retrieval.OWD* method), 22  
`__init__()` (*covid19\_inference.data\_retrieval.RKI* method), 18  
`__init__()` (*covid19\_inference.data\_retrieval.RKIsituationreports* method), 20  
`__init__()` (*covid19\_inference.data\_retrieval.retrieval.Retrieval* method), 23  
`_download_csv_from_source()` (*covid19\_inference.data\_retrieval.retrieval.Retrieval* method), 23  
`_fallback_handler()` (*covid19\_inference.data\_retrieval.retrieval.Retrieval* method), 23  
`_save_to_local()` (*covid19\_inference.data\_retrieval.retrieval.Retrieval* method), 24  
`_timestamp_local_old()` (*covid19\_inference.data\_retrieval.retrieval.Retrieval* method), 23  
`download_all_available_data()` (*covid19\_inference.data\_retrieval.GOOGLE* method), 21  
`download_all_available_data()` (*covid19\_inference.data\_retrieval.JHU* method), 16  
`download_all_available_data()` (*covid19\_inference.data\_retrieval.OWD* method), 22  
`download_all_available_data()` (*covid19\_inference.data\_retrieval.RKI* method), 18  
`download_all_available_data()` (*covid19\_inference.data\_retrieval.RKIsituationreports* method), 20

## F

`filter()` (*covid19\_inference.data\_retrieval.RKI* method), 19  
`filter_bundesland()` (*covid19\_inference.data\_retrieval.RKI* method), 20  
`filter_date()` (*covid19\_inference.data\_retrieval.JHU* method), 17

## B

`backup_instances()` (in module *covid19\_inference.data\_retrieval.retrieval*), 15

## C

`Cov19Model` (class in *covid19\_inference.model*), 7  
`covid19_inference` (module), 1  
`covid19_inference.data_retrieval.retrieval` (module), 15  
`covid19_inference.model` (module), 7  
`covid19_inference.plotting` (module), 25

## D

`delay_cases()` (in module *covid19\_inference.model*), 11

## G

`get_all_free_RVs_names()` (in module *covid19\_inference.plotting*), 25  
`get_changes()` (*covid19\_inference.data\_retrieval.GOOGLE* method), 21  
`get_new()` (*covid19\_inference.data\_retrieval.JHU* method), 17  
`get_new()` (*covid19\_inference.data\_retrieval.OWD* method), 22  
`get_new()` (*covid19\_inference.data\_retrieval.RKI* method), 19  
`get_possible_counties_states_regions()` (*covid19\_inference.data\_retrieval.GOOGLE* method), 22  
`get_possible_countries()` (*covid19\_inference.data\_retrieval.OWD* method), 22

`get_possible_countries_states()` (*covid19\_inference.data\_retrieval.JHU method*), 18  
`get_prior_distribution()` (*in module covid19\_inference.plotting*), 25  
`get_total()` (*covid19\_inference.data\_retrieval.JHU method*), 17  
`get_total()` (*covid19\_inference.data\_retrieval.OWD method*), 22  
`get_total()` (*covid19\_inference.data\_retrieval.RKI method*), 18  
`get_total_confirmed_deaths_recovered()` (*covid19\_inference.data\_retrieval.JHU method*), 16  
`GOOGLE` (*class in covid19\_inference.data\_retrieval*), 21

## H

`hierarchical_normal()` (*in module covid19\_inference.model*), 13

## J

`JHU` (*class in covid19\_inference.data\_retrieval*), 16

## L

`lambda_t_with_sigmoids()` (*in module covid19\_inference.model*), 13

## M

`make_change_point_RVs()` (*in module covid19\_inference.model*), 12  
`make_prior_I()` (*in module covid19\_inference.model*), 13  
`modelcontext()` (*in module covid19\_inference.model*), 8

## O

`OWD` (*class in covid19\_inference.data\_retrieval*), 22

## P

`plot_cases()` (*in module covid19\_inference.plotting*), 25  
`plot_hist()` (*in module covid19\_inference.plotting*), 25

## R

`Retrieval` (*class in covid19\_inference.data\_retrieval.retrieval*), 23  
`RKI` (*class in covid19\_inference.data\_retrieval*), 18  
`RKIsituationreports` (*class in covid19\_inference.data\_retrieval*), 20

## S

`SEIR()` (*in module covid19\_inference.model*), 9

`set_data_dir()` (*in module covid19\_inference.data\_retrieval.retrieval*), 15  
`SIR()` (*in module covid19\_inference.model*), 9  
`student_t_likelihood()` (*in module covid19\_inference.model*), 8

## W

`week_modulation()` (*in module covid19\_inference.model*), 12