

# deepregression: a Flexible Neural Network Framework for Semi-Structured Deep Distributional Regression

David Rügamer

LMU Munich

Ruolin Shen Christina Bukas Lisa Barros de Andrade e Sousa Dominik Thalmeier

Helmholtz AI Munich

Nadja Klein Chris Kolb  
Humboldt-Universität zu Berlin

Florian Pfisterer Bernd Bischl  
LMU Munich

Christian L. Müller

LMU Munich, ICB, Helmholtz Zentrum Munich,  
CCM, Flatiron Institute, New York

---

## Abstract

This paper describes the implementation of semi-structured deep distributional regression, a flexible framework to learn distributions based on a combination of additive regression models and deep neural networks. **deepregression** is implemented in both R and Python, using the deep learning libraries **TensorFlow** and **PyTorch**, respectively. The implementation consists of (1) a modular neural network building system for the combination of various statistical and deep learning approaches, (2) an orthogonalization cell to allow for an interpretable combination of different subnetworks as well as (3) pre-processing steps necessary to initialize such models. The software package allows to define models in a user-friendly manner using distribution definitions via a formula environment that is inspired by classical statistical model frameworks such as **mgcv**. The packages' modular design and functionality provides a unique resource for rapid and reproducible prototyping of complex statistical and deep learning models while simultaneously retaining the indispensable interpretability of classical statistical models.

*Keywords:* additive predictors; deep learning; effect decomposition; orthogonal complement; penalization and smoothing.

---

## 1. Introduction

In regression analysis, the overarching objective is to obtain information about the conditional distribution of a response variable given a set of explanatory variables. Yet, many regression approaches only focus on modeling the mean of this conditional distribution. Such approaches often cannot account for heteroscedasticity of the conditional distribution or its possible de-

pendence on features (covariates). As a result, methodological research has brought forward a number of distinct approaches, often referred to as “distribution, distributional or density regression”. The common ground of all such contributions is to allow characteristics of the conditional distribution beyond the mean to vary with (different types) of available features. Common examples include extensions of generalized additive models (GAMs; see, e.g., Wood 2017) by making all distributional parameters of a parametric density feature-dependent (e.g., generalized additive models for location, scale, and shape (GAMLSS); Rigby and Stasinopoulos 2005), but also approaches avoiding the parametric distributional assumptions such as Bayesian nonparametrics (Dunson 2010), mixtures of experts (Jacobs, Jordan, Nowlan, and Hinton 1991), distributional regression or conditional transformation models (Hothorn, Kneib, and Bühlmann 2014). A further branch focuses on directly relating local properties of the conditional distribution to features, indexed by the quantile or expectile level; see Koenker and Bassett (1978); Newey and Powell (1987) for the original references and, e.g., Koenker (2005); Schnabel and Eilers (2009) for more recent contributions in this area.

Despite being more flexible than mean regression approaches, many of these distributional regression methods are not scalable to either large number of observations or large number of features. The semi-structured deep distributional regression (SDDR) framework proposed by Rügamer, Kolb, and Klein (2020) and implemented in the package **deepregression** overcomes this limitation while allowing for the flexibility of distributional regression. Within the package it is possible to learn the distributional parameters of a wide range of parametric densities based on a combination of structured additive regression models and deep neural networks. **deepregression** thereby offers a scalable and flexible alternative to existing distributional regression approaches.

### 1.1. Semi-Structured deep distributional regression

SDDR follows the ideas of distributional regression based on parametric distributions as, e.g., done in a fully Bayesian framework by Klein, Kneib, Lang, and Sohn (2015), in order to estimate the entire conditional distribution of a continuous, discrete, mixed, or multivariate response variable  $\mathbf{Y}$ . This is done implicitly through modelling the corresponding distributional parameters. Apart from its application to transformation models (Baumann, Hothorn, and Rügamer 2020), SDDR has also been applied recently to survival regression (Kopper, Pölsterl, Wachinger, Bischl, Bender, and Rügamer 2021), mixture models (Rügamer, Pfisterer, and Bischl 2020), and to combine epidemiological models and graph neural networks (Fritz, Dorigatti, and Rügamer 2021).

To be more precise, SDDR constitutes a distributional learning framework where each distributional parameter  $\theta_k \equiv \theta_k(\eta_k), k = 1, \dots, K$  of an arbitrary pre-specified parametric distribution  $\mathcal{D}(\theta_1, \dots, \theta_K)$  can be modeled through (potentially different) additive predictors  $\eta_k$ .<sup>1</sup> Each of these predictors can be defined as a sum of structured predictors, which are either (penalized) linear effects or (penalized) smooth effects, or arbitrary additional (unstructured)

---

<sup>1</sup>It is assumed that  $\mathcal{D}$  has an absolute continuous density with existing finite first derivatives with respect to all distributional parameters.

neural network predictors of features  $\mathbf{x}$ ,  $\mathbf{z}$ ,  $\mathbf{u}$ :

$$\eta_k \equiv \eta_k(\mathbf{x}, \mathbf{z}, \mathbf{u}) = \mathbf{x}^\top \mathbf{w} + \sum_{j=1}^J f_{k,j}(\mathbf{z}) + \sum_{l=1}^L d_{k,l}(\mathbf{u}), \quad (1)$$

where the first part on the right-hand side are the linear effects of features  $\mathbf{x}$  and unknown weights  $\mathbf{w}$ , the second part are the smooth effects of one or several features  $\mathbf{z}$ , and the last represents the (deep) neural networks of one or several features  $\mathbf{u}$ . The smooth effects  $f_{k,j}$  are defined by appropriate basis function representations  $B_1, \dots, B_M$ , such as certain spline basis functions (see e.g. Wood 2017, for A recent overview). In summary, the following transformations are applied:

$$(\mathbf{x}, \mathbf{z}, \mathbf{u}) \xrightarrow{(1)} \eta_k(\mathbf{x}, \mathbf{z}, \mathbf{u}) \xrightarrow{h_k} \theta_k(\eta_k) \longrightarrow \mathcal{D}(\theta_1, \dots, \theta_k, \dots, \theta_K).$$

Here,  $h_k(\cdot)$  are one-to-one transformations mapping  $\eta_k$  from the real line to the respective parameter spaces of  $\theta_k$  (also known as response functions). For example, choosing  $h_k = \exp(\cdot)$  can be used to model a positive scale parameter. To ensure identifiability and interpretability in the additive predictors, the SDDR framework uses an orthogonalization cell that allows to separate effects in the structured part from the unstructured part of the predictors. This only requires the neural network predictors  $d_l$  to have a fully-connected layer as penultimate layer (cf. Figure 6).

The SDDR framework comprises many special cases such as linear and generalized linear models, GAMs, GAMLSS, and mixture models, but also allows to extend those classes via additional neural network predictors. In order to allow practitioners to make use of such a generic framework, **deepregression** implements the SDDR framework in an accessible and modular manner in both R and Python. Before going into the details of the corresponding packages, we first introduce our running example used throughout this tutorial.

## 1.2. Case study: Airbnb price listing data

As a running example, we will use price listings of apartments on Airbnb, available at <http://insideairbnb.com/get-the-data.html>. The data include various data modalities, from numeric variables such as latitude and longitude, to integer variables such as number of bedrooms, to textual information, including a room description, dates, and an image of each property. This makes the data set an ideal model example for the application of **deepregression**. In order to reduce computational complexity, we will focus on apartments in Munich and show exemplary code for demonstration purposes.

The cleaned `airbnb` Munich data consist of 3,504 observations and 74 features. In order to demonstrate the functionality of our package, we will focus on a handful of features, namely: `description` (a textual description of the property), `host_since` (a date variable describing the host's membership duration), `longitude` and `latitude`, `price` (the price of the apartment; our response variable), `review_score_value` (a review score), `host_response_time` (the time until the host response to inquiries; will also be used as an alternative response variable), and the properties' images (which are named according to the `id` variable in the dataset).

```
R> url <- "github.com/davidruegamer/airbnb/raw/main/munich_clean_text.RDS"
```

```
R> destfile <- file.path(getwd(), "munich.RDS")
R> download.file(url, destfile)
R> airbnb <- readRDS("munich.RDS")
R> airbnb$days_since_last_review <- as.numeric(
+   difftime(airbnb$date, airbnb$last_review)
+ )
```

The remainder of this paper is structured as follows: Section 2 details the core functions and building blocks in the corresponding R package. Section 3 provides insights into the Python-native implementation of the SDDR framework. A brief summary and conclusion is given in the final Section 4.

## 2. R package

**deepregression** is constructed in a Lego system fashion, allowing to specify a large number of distributions through the three different additive predictor types defined in (1). A three-parameter distribution, such as, e.g., the location, scale, and shape-parameterized  $t$ -distribution, can have different predictors defined for each of the  $K = 3$  parameters. Here, the scale and shape parameter could be nuisance parameters that are modeled using a deep neural network (DNN) of all features, while the location parameter is assumed to have specific (non-)linear effects in certain features and is thus modeled in a structured manner. Figure 1 visualizes three of the different building blocks in the Lego system. From a network perspective, the main building blocks are (i) a structured linear layer, (ii) a structured non-linear layer, (iii) an orthogonalization cell, (iv) distribution layer(s). The distribution layer (Figure 1(c)) defines the final distribution using one or more subnetworks for each parameter  $\theta_k$ . The network is then trained using the negative log-likelihood  $-\log \mathbf{p}_{\mathcal{D}(\boldsymbol{\theta})}(\mathbf{y})$ . In every subnetwork, the additive predictors  $\eta_k$  can be defined using (penalized) linear effects (Figure 1(a)) or (penalized) smooth effects (Figure 1(b)). These structured layers can be defined using penalties, such as  $L_1$ - or  $L_2$ -penalties for linear layers, and structured penalties induced by penalty matrices  $\mathbf{P}_j$ ,  $j = 1, \dots, J$  for the structured non-linear layer.

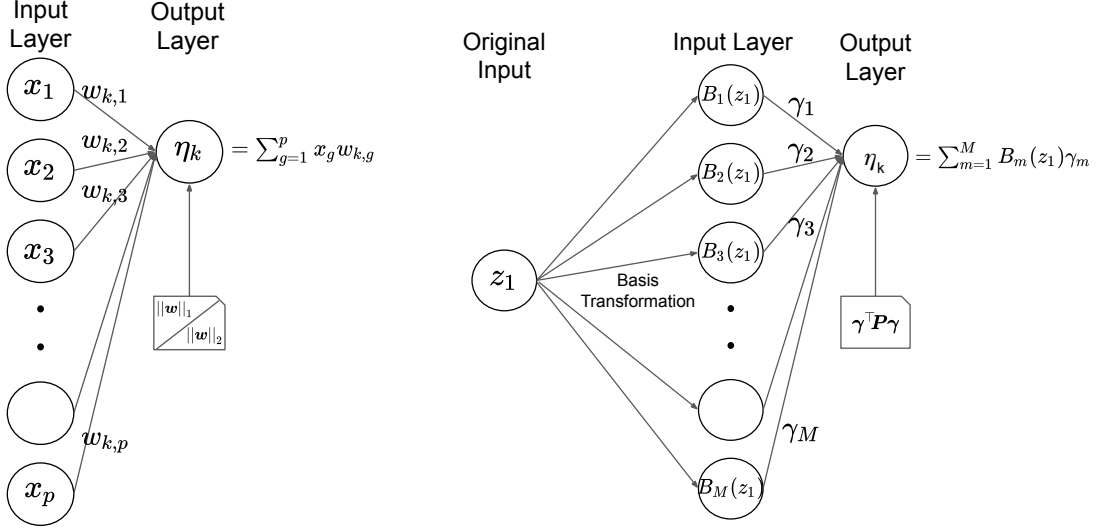
The network can also be turned into a Bayesian neural network by allowing each layer to have a distribution defined over its weights. These weights are used as prior distribution assumptions and an approximate Bayesian posterior is calculated using variational inference (Blundell, Cornebise, Kavukcuoglu, and Wierstra 2015; Tran, Mike, van der Wilk, and Hafner 2019). Furthermore, the distributional layers can be combined to define a richer distribution family, e.g., to create zero-inflated distributions or mixture of distributions.

While the structured linear layer can take inputs in their original form, the DNNs have to be user-specified to correctly process the inputs. The structured non-linear layer requires additional pre-processing and is described below as part of the three main core functionalities.

### 2.1. Core functions

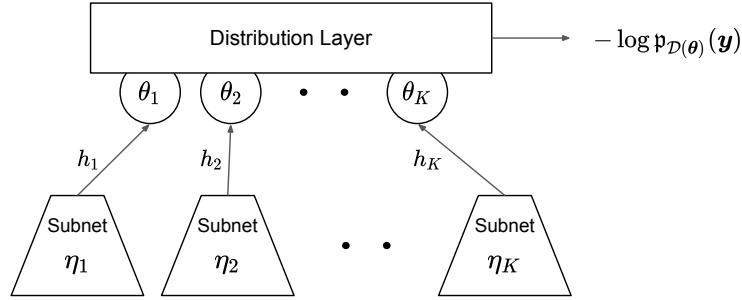
#### *Formula interface*

In the formula interface, each distribution parameter can be specified by a symbolic additive predictor. A model with two parameters  $\theta_1, \theta_2$  for location and scale of a distribution  $\mathcal{D}$  with



(a) A structured linear network: linear combination of input features  $x_j$  with optional  $L_1$ - or  $L_2$ -penalty on the weights  $w_{k,j}$

(b) A structured non-linear network: linear combination of (basis-transformed) features  $B_m(z_1)$  via  $M$  basis functions  $B_m, m = 1, \dots, M$  with optional matrix-based quadratic penalty  $\mathbf{P}$  on the weights  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_M)^\top$ .



(c) A distributional layer learned by  $K$  subnetworks: each parameter  $\theta_k$  of  $\mathcal{D}$  with density  $\mathfrak{p}_{\mathcal{D}}$  can be learned using (a), (b) and/or a DNN.

Figure 1: Illustration of network components that can be combined using **deepregression**. DNNs are left out as these can be arbitrarily specified and combined with additive predictors as long as the penultimate layer in the DNN corresponds to a fully-connected layer.

$\theta_1 = w_0 + w_1 x + f_{1,1}(z)$  and  $\theta_2 = f_{2,1}(u) + d_{2,1}(u)$  in R translates to:

```
R> list_of_formulae = list(
+   loc = ~ 1 + x + s(z, bs="tp"),
+   scale = ~ 0 + s(u, bs="ps") + dnn(u)
+ )
```

This specifies an intercept (also called bias term) for the location parameter as well as a linear effect for  $x$  and a thin-plate regression spline for  $z$ . The scale parameter is specified without

bias term but using another non-linear effect of  $u$  and a DNN predictor for inputs  $w$ , where `dnn` is a model defined separately with suitable input shape. In the network, this formula will create two subnetworks, both with one output unit (the location and scale parameter), which is then used to define the distribution.

### Case study: formula interface

We start with a simple additive model to predict the log-price

```
R> y = log(airbnb$price)
```

using a normal distribution parameterized by location and scale. We learn the location using an intercept and a tensor-product spline for `longitude` and `latitude`, while estimating the scale to be constant. The respective formula is

```
R> list_of_formulae = list(
+   loc = ~ 1 + te(latitude, longitude),
+   scale = ~ 1
+ )
```

In contrast to other regression software, **deepregression** also allows to specify (deep) neural networks in the formula. Details on the usage of the formula will be provided in Section 2.1.5.

### *Network initialization*

After model formula specification by the user, the SDDR model is initialized via the `deepregression` function. The required arguments are the response `y`, the data frame or list `data` including all features, the distribution `family`, the `list_of_formulae`, and the `list_of_deep_models` to define the (deep) neural networks that are used to model unstructured effects.

### Case study: network initialization

```
R> mod_simple <- deepregression(
+   y = y,
+   data = airbnb,
+   list_of_formulae = list_of_formulae,
+   list_of_deep_models = NULL
+ )
```

The result `mod` is of class `deepregression`, a list of length two, where the first element `model` is a **keras** model and the second element `init_params` is a list of objects required for model fitting, tuning, plotting, and other post-processing functionalities. The `model` element contains the references to the respective Python objects and is independent of the actual data. As the dimensions of the `model` can only be determined after evaluation of the model formula, yet **keras** models do not store any data, the second list item is used for convenience to allow the user the direct access to the pre-processed data (in the shape as required by the **keras** model) and to simplify calls to other methods such as `fit`.

```
R> mod_simple
```

```
Model
```

```
Model: "model"
```

```
-----
Layer (type)           Output Shape      Param #  Connected to
=====
input_1 (InputLayer)   [(None, 25)]      0
-----
input_2 (InputLayer)   [(None, 1)]       0
-----
structured_nonlinear_1 (D (None, 1)      25      input_1[0] [0]
-----
structured_linear_2 (Dens (None, 1)      1        input_2[0] [0]
-----
concatenate (Concatenate) (None, 2)      0        structured_nonlinear_1[0] [0]
                                         structured_linear_2[0] [0]
-----
distribution_lambda (Dist ((None, 1), (None 0      concatenate[0] [0]
=====
```

```
Total params: 26
```

```
Trainable params: 26
```

```
Non-trainable params: 0
```

```
-----
Model formulae:
```

```
-----
loc :
~1 + te(latitude, longitude)
scale :
~1
```

### *Pre-processing for structured non-linear layers*

Structured non-linear effects can be fitted in **deepregression** by defining a corresponding pre-processing function that generates a design matrix based on (user-defined) basis functions and a penalty matrix that is used to regularize the roughness of the estimated non-linear functions. In R, these objects are computed by calling the function `smooth.construct` from the package **mgcv** (Wood 2001).

### **Case study: Pre-processing**

In the above defined model `mod`, the smooth terms from **mgcv** are stored in the `init_params` element and are used for model fitting, prediction, and plotting.

```
R> str(mod$init_params$parsed_formulae_contents[[1]],1)
List of 3
```

```

$ linterms      :'data.frame': 3504 obs. of  1 variable:
$ smoothterms:List of 1
$ deeperterms   : NULL
- attr(*, "formula")=Class 'formula' language ~1 + te(latitude, longitude)
.. ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
- attr(*, "df")=List of 1
- attr(*, "variable_names")= chr [1:73] "id" "scrape_id" "last_scraped" ...
- attr(*, "network_names")= chr "d"
- attr(*, "intercept")= logi TRUE
- attr(*, "defaultSmoothing")=function (st, this_df)
- attr(*, "zero_cons")= logi TRUE

```

### *Specification of the family*

The `family` argument specifies the respective family to be learned. Possible choices can be searched for in `?dr_families`, comprising over 30 different options, including those from the exponential family (such as Gaussian, Bernoulli, Poisson, beta, gamma, and multinomial distribution), but also multivariate and mixture distributions.

**deepregression** can also define mixtures of families, as proposed by [Rügamer \*et al.\* \(2020\)](#), or custom distributions. Details can be found in the Section 2.6.

### *Formula contents*

The formulas defining each distributional parameter can be specified in the same way as in `mgcv::gam`, including `s`-terms, `ti`- and `te`-terms as well as further specifications of those smooth effects like the smoothing basis (see `?mgcv::s`). Factor variables in the formula are treated also the same way as in conventional regression packages by creating an encoded matrix (usually dummy-encoding / one-hot encoding). The exclusion of the intercept in one linear predictor can be defined as per usual using `0` in the formula: `~ 0 + ...` or via `-1`. An example using the `airbnb` data will be given in the next subsection.

### *Specification of the DNNs*

The DNNs specified in the `list_of_formulae` must also be passed as named list elements in the `list_of_deep_models`, i.e., each model term in the `list_of_formulae` that can not be recognized as intercept, linear effect of a `gam`-term must be listed in `list_of_deep_models`. The named list `list_of_deep_models` contains a list of DNNs, each specified as a function. These functions take as many inputs as defined for the neural network and end with a fully-connected layer with one hidden unit (or any other layer that results in the appropriate amount of outputs). The following code specifies an exemplary function for the later example using the pipe `%>%` notation:

```

R> deep_model <- function(x)
+   {
+     x %>%
+     layer_dense(units = 24, activation = "relu", use_bias = FALSE) %>%
+     layer_dropout(rate = 0.2) %>%

```



```

+     layer_dense(units = 12, activation = "relu") %>%
+     layer_dropout(rate = 0.2) %>%
+     layer_dense(units = 1, activation = "linear")
+   }

```

This defines a two hidden-layer DNN with dropouts between each layer. To ensure identifiability when structured regression terms and a DNN share some input features, an orthogonalization cell is automatically included before the last dense layer. In this case it is required to use a 'linear' activation function as in this example.

### Case study: formulas and DNNs

The following code shows an exemplary specification including all possible formula terms. We only define  $\theta_1$  using an intercept, categorical effects for `beds`, two splines (P-spline and thin-plate regression spline) for the number of guests and days since last review and a DNN for the review score and the number of reviews per month.

```

R> mod <- deepregression(
+   y = y,
+   data = data,
+   list_of_formulae = list(
+     location = ~ 1 + beds + s(accommodates, bs = "ps") +
+       s(days_since_last_review, bs = "tp") +
+     deep(review_scores_rating, reviews_per_month),
+     scale = ~1),
+   list_of_deep_models = list(deep = deep_model)
+ )

```

This model definition translates to an additive regression model where the location of the model is defined by an intercept 1, a linear effect of the factor variable `beds`, a P-spline for variable `accommodates`, a thin-plate regression spline for variable `days_since_last_review` and a DNN for variables `review_scores_rating`, `reviews_per_month`.

## 2.2. Model fitting and tuning

### *Model fitting*

Once the model `mod` has been set up, the neural network can be trained using the `fit` function. The `fit` function is a wrapper of the corresponding `fit` function for `keras` models and inherits the `keras::fit` arguments. More specifically,

- `epochs` to specify the number of iterations,
- `callbacks` to specify information that is called after each epoch or batch (used, e.g., for early stopping),
- `validation_split` to specify the amount of data (in  $[0,1)$ ) that is used to validate the model (while the rest is used to train the model),

- `validation_data` to optionally specify any predefined validation data,
- `verbose` to define if progress is printed in the console and
- `view_metrics` to activate a interactive plot window when using RStudio.

Several other arguments are available, such as a logical argument `early_stopping` to activate early stopping and `patience` to define the patience used in early stopping.

### Case study: model fitting

```
R> mod %>% fit(
+   epochs = 100,
+   verbose = FALSE,
+   view_metrics = FALSE,
+   validation_split = 0.2
+ )
```

To inspect the fitted model, we can either use `coef` or `plot`, as demonstrated previously, or extract the fitted values of the model:

```
R> fitted_vals <- mod %>% fitted()
R> cor(fitted_vals, y)
      [,1]
[1,] 0.5805392
```

### *Model tuning*

In a similar manner as `fit`, **deepregression** offers a cross-validation function `cv` for model tuning. This can be used to fine-tune the model by, e.g., changing the formula(s), changing the DNN structure or defining the amount of smoothness (using the `df` argument in `deepregression`; see the next section for details). Folds in the `cv` function can be specified using the `cv_folds` argument. This argument takes either an integer for the number of folds or a list where each list item is again a list of two, one element with data indices for training and one with indices for testing. The `patience` for early stopping can be set using the respective argument.

### Case study: cross-validation

As an example we use 100 iterations and three folds. The results of the cross-validation can be plotted using `plot_cv`.

```
R> res_cv <- mod %>% cv(
+   plot = FALSE,
+   cv_folds = 3,
+   epochs = 100
+ )
R> plot_cv(res_cv)
```

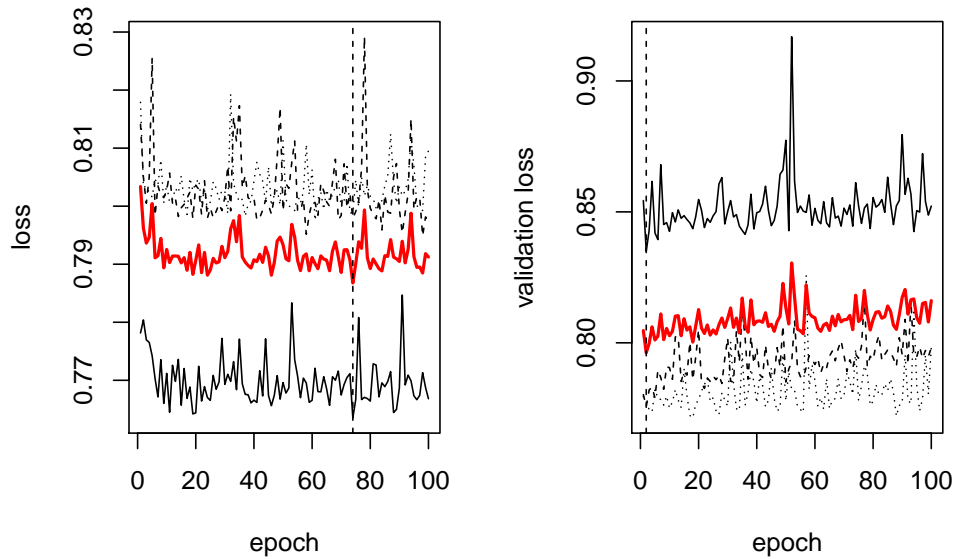


Figure 2: Cross-validation result. Horizontal black lines indicate the different folds over different epochs (x-axis) and the corresponding loss (left) or validation loss (right). The red line shows the mean loss and validation loss over all iterations. The dashed vertical line indicates the epochs with the lowest average loss and validation loss in the respective plots.

### 2.3. Methods overview

Apart from `fit` and `cv`, there are additional methods that can be applied to a `deepregression` object:

- `coef`: The `coef` function extracts the coefficients (network weights) of all layers with structured definition, i.e., coefficients of linear or additive effects. Using the argument `type`, the user can specify what type of coefficients to return (possible choices: "linear" for linear effects and "smooth" for basis coefficients of the specified smooth terms) and using the argument `params` allows to select which of the distribution parameter's coefficient the user wants to return.
- `plot`: The `plot` function can be applied to `deepregression` objects to plot the estimates of non-linear effects, i.e., splines and tensor product splines. Using `which` a specific effect can be selected using the corresponding integer in the structured part of the formula (if `NULL` all effects are plotted), while the integer given for `which_param` indicates which distribution parameter is chosen for plotting. Via the argument `plot = FALSE`, the function can also be used to just return the data used for plotting.
- `predict`: The `predict` function works as for other (`keras`) models. It either just takes the model as input and returns the predictions (per default the distribution's expectation) for observations in the training data set, or, when supplied with new data, produces the corresponding predictions for the new data. As `deepregression` learns a distribution and not only a mean prediction, the user can choose via the argument `apply_fun` what type of distribution characteristic is to be predicted (per default `apply_fun = tfd_mean` predicts the mean of the distribution).

- **mean**, **sd** and **quantile**: The functions are convenience functions that directly return the mean, standard deviation, or certain quantiles of the learned distribution. All three functions work for the given training data as well as for newly provided **data**. The **quantile** function has an additional argument **probs** to provide the quantile(s) of interest.
- **get\_distribution**: Instead of returning summary statistics, **get\_distribution** returns the whole **TensorFlow** distribution with all its functionalities, such as sampling random numbers from the fitted distribution, computing the PDF or CDF, etc.
- **log\_score**: The **log\_score** function directly returns the evaluated log-likelihood based on the estimated parameters and the provided data (including a response vector **this\_y**). If **data** and **this\_y** are not provided, the function will calculate the score on the training data.

### Case study: convenience functions

We first inspect the estimated linear effects.

```
R> coef(mod, type="linear")
$location
(Intercept)      beds.1      beds.2      beds.3      beds.4      beds.5
1.975234151  0.008926539  0.070396960  0.138633251  0.151451677  0.582612634
      beds.6      beds.7      beds.8      beds.9      beds.10      beds.11
0.221924141  0.175573707 -0.311866581 -0.123485431  0.524103165  1.533515930
      beds.12      beds.16      beds.18      beds.20      beds.30
-0.313976169 -0.611083746  1.997076750 -0.692837417 -0.524448991

$scale
(Intercept)
-0.6014107
```

For non-linear effects, the estimated smooth functions can be plotted using **plot** with result visualized in Figure 3.

While most of the other convenience functions are wrappers that access the underlying **keras** model and convert the quantity of interest to an R object, **get\_distribution** returns the fitted distribution, a Python **tfp.distributions**.

```
R> dist <- mod %>% get_distribution()
R> str(dist, 1)
```

Instead of working with the original distribution object, the user can work with convenience functions to get the quantity of interests directly, as demonstrated below.

```
R> meanval <- mod %>% mean()
R> q05 <- mod %>% quantile(probs = 0.05)
R> q95 <- mod %>% quantile(probs = 0.95)
```

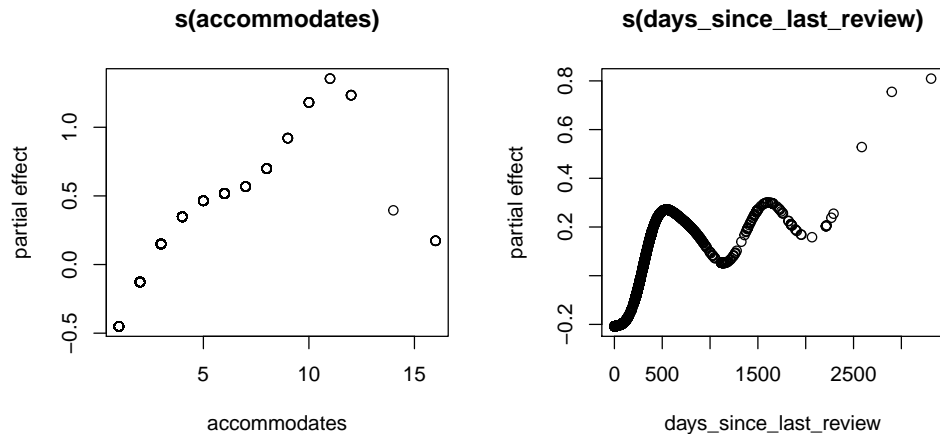


Figure 3: Case Study. Estimated partial effects of `accommodates` and `days_since_last_review` on the average log-price of an apartment.

```
R> plot(y, meanval, xlab = "True value", ylab = "Estimated mean value")
R> points(y, q05, pch="-", col="blue")
R> points(y, q95, pch="-", col="blue")
R> abline(0, 1, col = "red")
```

## 2.4. Penalties

**deepregression** allows for different penalties, including  $L_1$ -,  $L_2$ -, and other smoothing penalties as provided by the **mgcv** package. While the latter is implicitly created when using `s`-, `ti`- or `te`-terms in the formula, the  $L_1$ - and  $L_2$ -penalties are used to penalize linear predictor models without smooth terms by defining the amount of penalization via `lambda_lasso` or `lambda_ridge`, respectively. If both terms are used, this corresponds to an elastic net-type penalty (Zou and Hastie 2005). Since the model object returned by **deepregression** is a list where the first element is a **keras** model, additional penalties can always be specified after the model has been initialized. This is implemented by the `additional_penalty` argument which requires a function with one (unused) dummy argument and defines the penalty based on the **keras** model's `trainable_weights`.

### Case study: additional penalties

To access the model's weights, we can access the **keras** object in the previously created **deepregression** model `mod_simple`, use `mod_simple$model$trainable_weight` and, e.g., penalize all coefficients (weights) of the defined additive model using an  $L_1$ -penalty:

```
R> lambda <- 0.5
R> l1part <- tf$reduce_sum(tf$abs(mod_simple$model$trainable_weights[[1]]))
R> addpen <- function(x) lambda * l1part
R> mod_simple_pen <- deepregression(
+   y = y,
+   data = airbnb,
+   list_of_formulae = list_of_formulae,
```

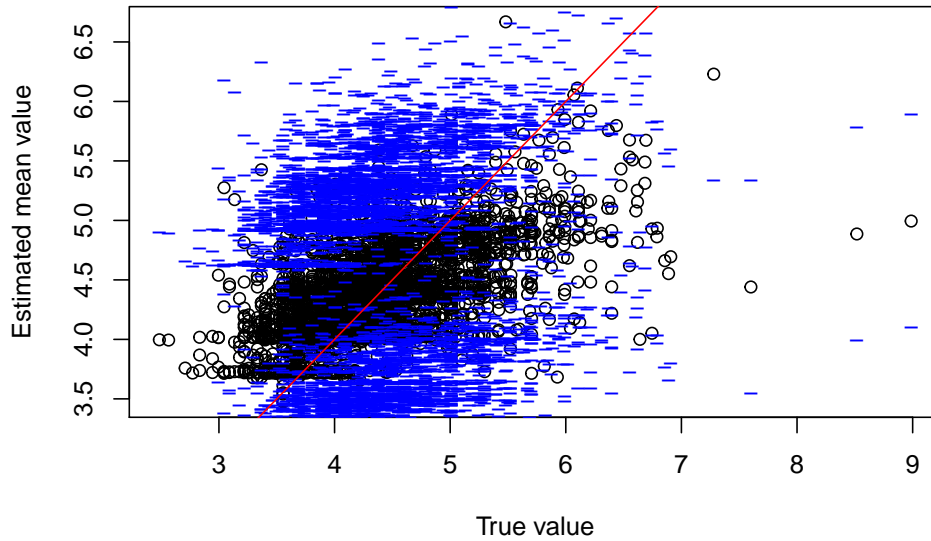


Figure 4: Visualization of estimated means and quantiles (black and blue points, respectively) against true values (x-axis).

```
+ list_of_deep_models = NULL,
+ additional_penalty = addpen
+ )
```

Note that the internal function must return a scalar value, and mathematical operations should be specified using **TensorFlow** functions, available in `tf` (and loaded together with **deepregression**).

### *Smoothing penalties*

Apart from specifying the smoothness of **mgcv**'s smooth terms in the formulas manually (see `mgcv::s` for more details), three further options are available in **deepregression**. The first option is to use a single degrees-of-freedom specification using the argument `df` in **deepregression**. Using Demmler-Reinsch orthogonalization (see, e.g., [Ruppert, Wand, and Carroll 2003](#)), all smoothing parameters are then calculated based on this specification. This ensures that no smooth term has more flexibility than the other term. When `df` is left unspecified, the default `df` is the smallest basis dimension among all smooths (leaving the least flexible smooth unpenalized, while the others are penalized to have the same degrees-of-freedom as this one). The second option in **deepregression** is to specify a list of length `length(list_of_formulae)`. This not only allows to specify different `df` for each distribution parameter, but also to specify different `df` for each smooth term in each formula by providing a vector of the same length as the number of smooths in the parameter's formula. The third option is based on a custom optimizer designed for additive model fitting in batch mode. This optimizer estimates the smooths directly and can be defined by using the argument `fsbatch_optimizer = TRUE`.

**Details** The definition of the degrees-of-freedom can be changed using the `hat1` argument. When set to `TRUE`, the `df` are assumed to be the sum of the diagonal terms of the hat matrix

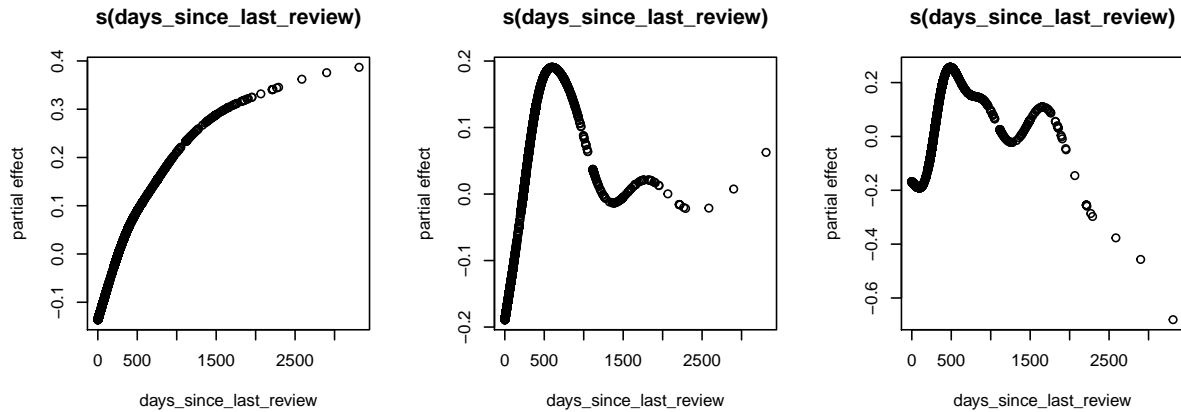


Figure 5: Case study. Comparison of different  $df$  settings for a single smooth term model.

$\mathbf{H}$  of the corresponding smooth term. The default is `FALSE`, yielding the more common definition of the effective degrees-of-freedom  $df = \text{trace}(2\mathbf{H} - \mathbf{H}\mathbf{H})$ . As the penalty in **TensorFlow** is added to each single observation, the penalty term is per default divided by the number observations, preventing too strong penalization of the log-likelihood. In certain situations, it might be necessary to scale the penalty differently. This can be done using the argument `sp_scale`, which is  $1/\text{NROW}(y)$  by default to account for batch training. If training is done in full batch mode (batch size is equal to the total number of observations), we recommend setting `sp_scale` to 1.

### Case study: smoothing penalties

In the following we use a single smooth term to demonstrate the effect of  $df$ .

```
R> forms <- list(loc = ~ 1 + s(days_since_last_review), scale = ~ 1)
R> args <- list(y = y, data = airbnb, list_of_formulae = forms,
+   list_of_deep_models = NULL)

R> mod_df_low <- do.call("deepregression", c(args, list(df = 3)))
R> mod_df_med <- do.call("deepregression", c(args, list(df = 6)))
R> mod_df_max <- do.call("deepregression", c(args, list(df = 10)))

R> mod_df_low %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)
R> mod_df_med %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)
R> mod_df_max %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)

R> par(mfrow=c(1,3))
R> plot(mod_df_low)
R> plot(mod_df_med)
R> plot(mod_df_max)
```

## 2.5. Neural network settings

Since **deepregression** constitutes a holistic neural network, certain settings and advantages

from deep learning can be made use of.

### *Shared DNN*

In addition to the DNN specifications, introduced in the Section 2.1.6, **deepregression** allows to share one DNN between some or all distribution parameters. This reduces the number of parameters to be estimated. The following example will demonstrate how to share a DNN (an LSTM model) between two parameters. The number of output units has to be equal to the number of parameters learned in this case. Sharing a network requires its inclusion in the `list_of_deep_models` and an additional list with formulas for the argument `train_together`. Analogously to `list_of_formulae`, this list provides the information which networks are used for which distribution parameters.

#### **Case study: Shared DNN**

In the following example, we will learn both location and scale parameter using an embedding of words contained in the room description. For text preparation of the original data, we refer the reader to the Appendix. We use the already tokenized text contained in the matrix `texts` in the `airbnb` data. We now define our neural network. The network consists of an embedding layer that processes the 1000 different tokenized words and learns a vector representation of the words in a 100-dimensional space. This representation is then reduced in the second dimension using the mean. Following this operation in the so-called lambda-Layer, we use two fully-connected layers with 20 and two units and dropout in between. The two units in the last layer correspond to the two parameters that are supposed to be learned from the text information.

```
R> embd_mod <- function(x) x %>%
+   layer_embedding(input_dim = 1000,
+     output_dim = 100) %>%
+   layer_lambda(f = function(x) k_mean(x, axis = 2)) %>%
+   layer_dense(20, activation = "tanh") %>%
+   layer_dropout(0.3) %>%
+   layer_dense(2)
```

In addition to the DNN, we will only use intercepts in both distribution parameters. Hence, the formulas we provide to the `list_of_formulae` is the following:

```
R> form_lists <- list(
+   location = ~ 1,
+   scale = ~ 1
+ )
```

To not copy and paste the `embd_mod` into every parameter formula, we create another list which is then provided as value for `train_together`:

```
R> embd_list <- list( ~ embd_mod(texts))[rep(1,2)]
```



```
R> mod <- deepregression(
+   y = y,
+   list_of_formulae = form_lists,
+   list_of_deep_models = list(embd_mod = embd_mod),
+   train_together = embd_list,
+   data = airbnb
+ )
```

This example creates an additive model where both the mean and the scale parameter are trained by the same network based on the text information `texts` and in the last layer, the latent features extracted from those descriptions are then split and added to the respective additive predictor. Alternatively, we could also define the network with a single unit output and add it directly to either one or both parameters:

```
R> embd_mod <- function(x) x %>%
+   layer_embedding(input_dim = 1000,
+     output_dim = 100) %>%
+   layer_lambda(f = function(x) k_mean(x, axis = 2)) %>%
+   layer_dense(20, activation = "tanh") %>%
+   layer_dropout(0.3) %>%
+   layer_dense(1)
```

```
R> form_lists <- list(
+   location = ~ 1 + embd_mod(texts),
+   scale = ~ 1
+ )
```

```
R> mod <- deepregression(
+   y = y,
+   list_of_formulae = form_lists,
+   list_of_deep_models = list(embd_mod = embd_mod),
+   data = airbnb
+ )
```

### *Optimizer and learning rate*

`deepregression` directly passes the `optimizer` to `keras`. It is therefore possible to specify any optimizer in `deepregression` that is available in `keras`, including `optimizer_adadelata`, `optimizer_adam`, `optimizer_rmsprop`, or standard stochastic gradient descent `optimizer_sgd`. The software uses Adam as default. The learning rate of the optimizer can be changed using the `learning_rate` argument, or, if another optimizer is defined, directly by using the respective argument of the above mentioned functions. For example, using Adadelata with learning rate 3 and decay 0.1 can be defined by passing the argument `optimizer = optimizer_adadelata(lr = 3, decay = 0.1)` to `deepregression`. This also overwrites the default `learning_rate` argument of `deepregression`.

### *Orthogonalization*

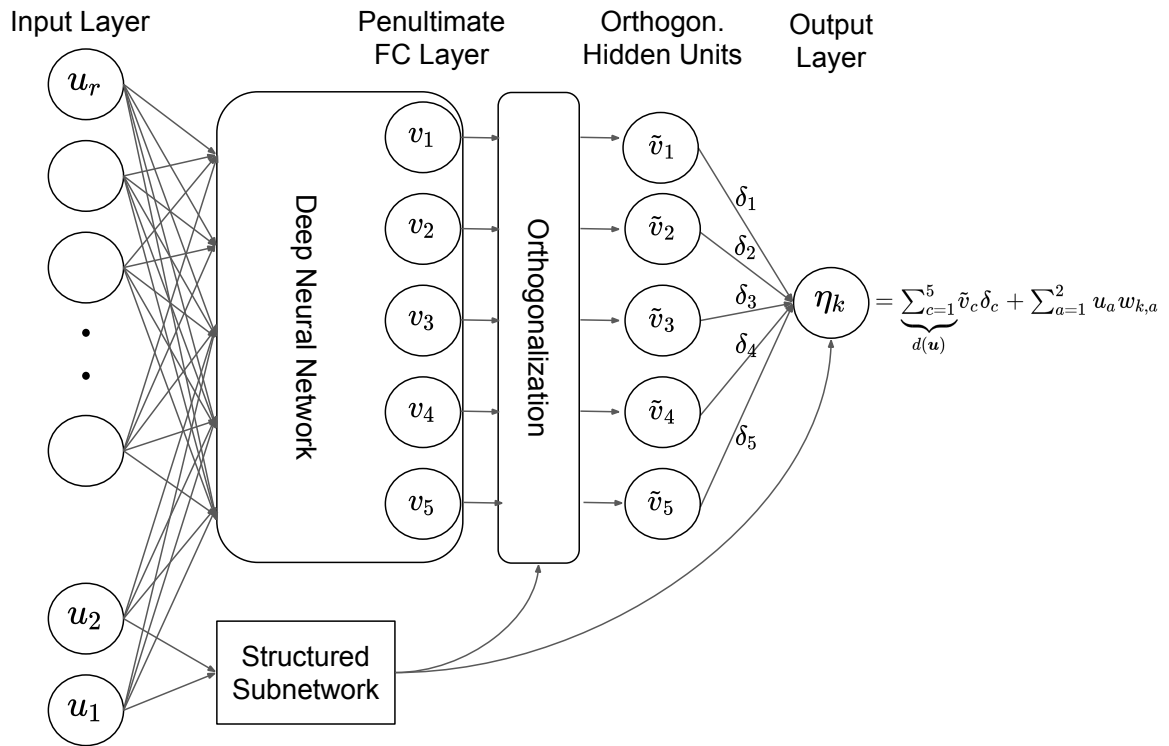


Figure 6: Orthogonalization cell: effects of partially shared features are learned in a DNN and in a structured subnetwork. The resulting latent features  $v_j$  in the DNN are then passed through an orthogonalization cell to ensure identifiability in the predictor  $\eta_k$ .

`deepregression` per default (`orthogonalize = TRUE`) orthogonalizes the DNNs in predictor formulas if these contain the same features as present in a structured partial effect in the same additive predictor (see Figure 6 for a schematic visualization).

### Case study: Orthogonalization

In the following, we introduce a toy example to demonstrate the orthogonalization property of SDDR. We define our data generating process as  $Y = 2x + \epsilon, \epsilon \sim \mathcal{N}(0, 1)$ :

```
R> toyX = rnorm(100)
R> toyY = 2*toyX + rnorm(100)
```

We here want to recover the linear effect of 2 of `toyX` in the presence of a DNN. We define this DNN as follows:

```
R> deep_model <- function(x)
+   {
+     x %>%
+       layer_dense(units = 24, activation = "relu", use_bias = FALSE) %>%
+       layer_dropout(rate = 0.2) %>%
+       layer_dense(units = 12, activation = "relu") %>%
+       layer_dropout(rate = 0.2) %>%
+       layer_dense(units = 1, activation = "linear")
+   }
```

Next, we fit a linear regression with linear effect for `toyX`. We pass the variable also to the `deep_model`, one time with orthogonalization and one time without orthogonalization.

```
R> forms <- list(loc = ~ -1 + toyX + deep_model(toyX), scale = ~ 1)
R> args <- list(
+   y = toyY,
+   data = data.frame(toyX = toyX),
+   list_of_formulae = forms,
+   list_of_deep_models = list(deep_model = deep_model)
+ )

R> mod_w_oz <- do.call("deepregression", c(args, list(orthogonalize = TRUE)))
R> mod_wo_oz <- do.call("deepregression", c(args, list(orthogonalize = FALSE)))

R> mod_w_oz %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)
R> mod_wo_oz %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)

R> cbind(
+   with = coef(mod_w_oz, params = 1)[[1]],
+   without = coef(mod_wo_oz, params = 1)[[1]],
+   linmod = coef(lm(toyY ~ 0 + toyX))
+ )
```

```

      with without linmod
toyX 2.198575 -0.87707 2.265347

```

In the model with orthogonalization, the `deep_model` is orthogonalized w.r.t. `toyX` to ensure identifiability of the structured linear term. We thus can recover the linear effect in `mod_w_oz` despite the presence of the DNN which could have captured a linear (or more complex) effects as well. By default, orthogonalization automatically extracts the terms that overlap in the DNNs and the structured model formula. For expert use, there is also a custom orthogonalization available (see Section 2.8)

## 2.6. Advanced usage

**deepregression** allows for several advanced model specifications and user inputs, briefly described next.

### *Offsets*

Several statistical models require an offset to be used in one or more linear predictors. These can be specified using the argument `offset` as a list of column vectors (i.e., matrices with one column) or `NULL` to include no offset. By default, no offset is used in any parameter.

### *Constraints*

For smooth effects, two options are currently available to constrain their estimation. These are inherited by **mgcv**. `absorb_cons` (default `FALSE`) will absorb identifiability constraints into the smoothing basis (see `?mgcv::smoothCon`) and `zero_constraint_for_smooths` will constrain all smooth to sum to zero over their domain, which is usually recommended to prevent identifiability issues (default `TRUE`).

### *Custom distribution function*

It is also possible to define custom distribution functions to be learned in **deepregression** using the `dist_fun` argument. To specify a custom distribution, define `dist_fun` as follows:

```

function(x){
  do.call(your_tfd_dist,
    lapply(1:ncol(x)[[1]],
      function(i) your_trafo_list_on_inputs[[i]](
        x[,i,drop=FALSE])
    )
  )
}

```

where `your_tfd_dist` is a distribution from **tfprobability** and `your_trafo_list_on_inputs` applies a transformation, e.g., `tf$exp()` to ensure that the parameter is defined on the correct domain.

## 2.7. Mixture models

**deepregression** can be used to fit mixtures of distributional regression models as proposed by Rügamer *et al.* (2020). In order to define a mixture, **deepregression** provides the function `mix_dist_maker` to create the mixture distribution in **TensorFlow Probability**, as required by the `dist_fun` argument of **deepregression**. The `list_of_formulae` for mixtures is composed of one entry for the mixing probability (first list entry), which can also depend on data, followed by the additive predictors for each parameter in each mixture. Here, the user has to specify first the formulas for all parameters of the first mixture component, then the second mixture component, and so on. Apart from the additional specification of the `dist_fun` argument, the argument `mixture_dist` must specify an integer for the number of mixtures. All other workflows remain the same.

### Case study: Mixture models

For demonstration, we take a model from before and make its distribution assumption slightly more flexible by defining a mixture of three normal distributions instead of a single Gaussian distribution. We first generate the respective mixture distribution:

```
R> mixdist <- mix_dist_maker(
+   nr_comps = 3,
+   dist = tfd_normal,
+   trafos_each_param = list(
+     function(x) x,
+     function(x) 1e-8 + tf$math$exp(x)
+   )
+ )
```

The function `mix_dist_maker` requires three arguments. The first one is the number of mixtures, the second one is the **TensorFlow Probability** distribution for which a mixture should be created, and the third one defines the response function  $h_k$  that are applied to the different parameters in one mixture component. Here, we use the identity function for the mean and an `exp(.)` response function for the standard deviation, plus a small constant to avoid numerical instabilities.

Having defined the mixture distribution, we take the first list of formulas used in this tutorial and extend it for the given mixture of three distributions:

```
R> list_of_formulae = list(
+   loc = ~ 1 + te(latitude, longitude),
+   scale = ~ 1
+ )

R> form_mix <-
+   c(
+     list(~ 1),
+     list_of_formulae[rep(1:2,3)]
+   )
```

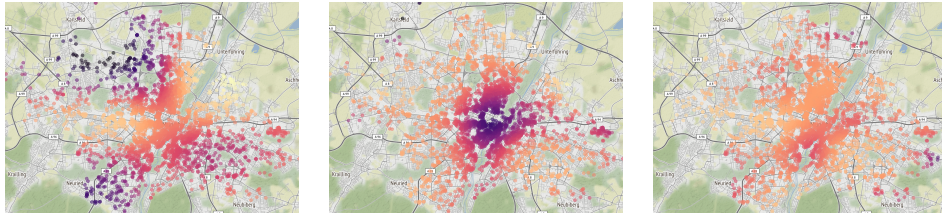


Figure 7: Case study. Effect of `longitude` and `latitude` on the average location price in the three estimated mixtures. Points correspond to apartments, black to purple points have an increasing effect on the average price while yellow to white points decrease the average price.

The `form_mix` now defines a constant mixture probability through the first list element and three equal normal distribution predictors, as defined in the original model with only one normal distribution.

```
R> mod_mix <- deepregression(
+   y = y,
+   list_of_formulae = form_mix,
+   list_of_deep_models = NULL,
+   data = airbnb,
+   mixture_dist = 3,
+   dist_fun = mixdist
+ )
```

Next, we run the model for a few iterations and look at the estimated effects:

```
R> mod_mix %>% fit(epochs=100, verbose=F)
R> par(mfrow=c(1,3))
R> do.call("grid.arrange", c(lapply(c(2, 4, 6), function(i)
+   map_plot(mod_mix, as.data.frame(airbnb), i)), list(nrow=1)))
```

## 2.8. Custom orthogonalization

If there is reason to orthogonalize a DNN w.r.t. a model term that is not explicitly present, the `%OZ%`-operator can be used. On the left side of the `%OZ%`-operator in the formula, the DNN must be given. On the right side of the operator, either a single model term (such as `x`, `s(x)`, `te(x,z)`) or a combination of model terms using brackets around all terms and separated with a `+` between each term must be supplied (e.g., `deep_model %OZ% (x + s(z))`). These terms must be structured terms.

### Case study: Custom orthogonalization

We use the previous toy example and pretend that there is another feature that can potentially influence the linear effect of `toyX` with the same information content but another name (so that it is not recognized automatically as a candidate for orthogonalization).

```
R> toyXinDisguise <- toyX
```

```

R> form_known <- list(loc = ~ -1 + toyX + deep_model(toyX), scale = ~ 1)
R> form_unknown <- list(
+   loc = ~ -1 + toyX + deep_model(toyXinDisguise),
+   scale = ~ 1
+ )
R> form_manual <- list(
+   loc = ~ -1 + toyX +
+     deep_model(toyXinDisguise) %OZ% (toyXinDisguise),
+   scale = ~ 1
+ )
R> args <- list(
+   y = toyY,
+   data = data.frame(toyX = toyX, toyXinDisguise = toyXinDisguise),
+   list_of_deep_models = list(deep_model = deep_model)
+ )

R> mod_known <- do.call(
+   "deepregression", c(args, list(list_of_formulae = form_known))
+ )
R> mod_unknown <- do.call(
+   "deepregression", c(args, list(list_of_formulae = form_unknown))
+ )
R> mod_manual <- do.call(
+   "deepregression", c(args, list(list_of_formulae = form_manual))
+ )

R> mod_known %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)
R> mod_unknown %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)
R> mod_manual %>% fit(epochs = 1000, early_stopping = TRUE, verbose = FALSE)

R> cbind(
+   known = coef(mod_known, params = 1)[[1]],
+   unknown = coef(mod_unknown, params = 1)[[1]],
+   manual = coef(mod_manual, params = 1)[[1]]
+ )

```

## 2.9. Working with images

The DNNs in SDDR can also incorporate images into the additive predictor of one or more parameters. Due to their size, images are usually not loaded into memory completely, but read in mini-batches from disk during training. For this purpose, **deepregression** offers a generator that creates small batches of images and other tabular information for training and prediction. The absolute path to the images must be given as string in the `data`. In addition, the user must specify the image size using the argument `image_var`, which is a named list of elements for each image source.

**Case study: Working with images**

The pictures of each of the apartments in the `airbnb` data can be downloaded from [github.com/davidruegamer/airbnb/tree/main/data/pictures](https://github.com/davidruegamer/airbnb/tree/main/data/pictures). Here, the folder 32 is the number of the city in the larger Airbnb data set, referred to in the introduction. Assuming the pictures are stored in the home directory in a dedicated folder called `airbnb`, we first create the absolute path as follows:

```
R> airbnb$image <- paste0("/home/user/airbnb/data/pictures/32/",
+   airbnb$id, ".jpg")
```

Next, we define an appropriate DNN architecture to convert images into latent features. Here, we use a convolutional neural network (CNN). First, we define a CNN block:

```
R> cnn_block <- function(
+   filters,
+   kernel_size,
+   pool_size,
+   rate,
+   input_shape = NULL
+ ){
+   function(x){
+     x %>%
+     layer_conv_2d(filters, kernel_size,
+       padding="same", input_shape = input_shape) %>%
+     layer_activation(activation = "relu") %>%
+     layer_batch_normalization() %>%
+     layer_max_pooling_2d(pool_size = pool_size) %>%
+     layer_dropout(rate = rate)
+   }
+ }
```

We then create the DNN function required by **deepregression** using a single block:

```
R> cnn <- cnn_block(
+   filters = 16,
+   kernel_size = c(3,3),
+   pool_size = c(3,3),
+   rate = 0.25,
+   shape(200, 200, 3)
+ )
R> deep_model_cnn <- function(x){
+   x %>%
+   cnn() %>%
+   layer_flatten() %>%
+   layer_dense(32) %>%
+   layer_activation(activation = "relu") %>%
```



```

+   layer_batch_normalization() %>%
+   layer_dropout(rate = 0.5) %>%
+   layer_dense(1)
+ }

```

Note that we have to provide the shape of the image in the first layer of the CNN, which is  $200 \times 200$  with three colour channels. Finally, we define an exemplary `deepregression` model using an additive predictor for the mean with both structured effects for, e.g., the room type as linear effect, and an unstructured effect by including the image information:

```

R> mod_cnn <- deepregression(
+   y = y,
+   list_of_formulae = list(
+     ~1 + room_type + bedrooms + beds +
+       deep_model_cnn(image),
+     ~1 + room_type),
+   data = airbnb,
+   image_var = list(image = list(c(200,200,3))),
+   list_of_deep_models = list(deep_model_cnn = deep_model_cnn),
+   optimizer = optimizer_adam(lr = 0.0001)
+ )

```

The training using images is usually slower, but also requires less iterations. For the given example, we can, e.g., use the following setting for model fitting:

```

R> mod_cnn %>% fit(
+   epochs = 100,
+   early_stopping = TRUE,
+   patience = 5,
+   verbose = FALSE
+ )

```

Checking the influence of the image in additive predictors is not straightforward. One strategy is to set all other tabular variables to zero or to the reference category and to qualitatively inspect predictions of the SDDR network. Note that, on small datasets, such as, e.g., the `airbnb` data for Munich considered here, SDDR networks with images have a tendency to overfit. To compare results of `deepregression` with its Python pendant `PySDDR` in the following, we print the estimated coefficients of tabular features

### 3. Python package

We also provide a Python-native implementation of the SDDR framework using `PyTorch`, called `PySDDR`. The core component of the package is the `SddrNet` class, which builds a dynamic network architecture. The Python package focuses on defining a class system that allows users familiar with `PyTorch` to easily extend the framework for their specific purposes. Modeling SDDR networks in Python directly can have various computational benefits (e.g.,

no communication needed between R and Python), and **PyTorch** itself can have benefits in comparison to **TensorFlow**. Prior to showing the application of the package in practice, we will provide a short overview of the module, its interface, and its functionality.

### 3.1. Overview of module

**SddrNet** offers a dynamic network architecture that depends on the user's inputs, i.e., the specified model distributions and formulas of the distributional parameters. The user needs to define a formula for each distributional parameter. Based on each of these formulas **SddrNet** builds parallel subnetworks, instances of the **SddrFormulaNet** class, one for each distributional parameter. The output of each **SddrFormulaNet** is a predicted additive predictor  $\eta_k$ . These predictors  $\eta_k, k = 1, \dots, K$  are collected by **SddrNet**, transformed based on the distribution's response functions  $h_k$ , and then passed to the distributional layer. From the latter, the penalized log-likelihood is computed, which is then back-propagated during training. A process diagram of the **SddrNet** module is available at <https://github.com/HelmholtzAI-Consultants-Munich/PySDDR>.

### 3.2. Interface

The user interacts with the package through the **Sddr** class. An overview of this class and its interaction with the rest of the package is available at <https://github.com/HelmholtzAI-Consultants-Munich/PySDDR>.

There are a number of inputs which need to be defined before training or testing can be performed. The user can either directly define these in a script or supply all the parameters via a config file. While there are various options that can be defined by the user, the most important input is a dictionary of formulas, one for each distribution parameter, analogous to the R package. **PySDDR** uses **Patsy** (<https://doi.org/10.5281/zenodo.592075>) in the backend to parse these formulas. Each formula is given as a string and follows Patsy's formula conventions. For example, a normal distribution with linear effect for **x1**, B-splines for further features and DNNs can be defined as

```
formulas = {
  'loc': '~ 1 + x1 + spline(x1, bs="bs", df=4) + d1(x3) + d2(x4, x5)',
  'scale': '~ -1 + spline(x3, bs="bs",df=4) + spline(x4, bs="bs",df=4)'
}
```

In contrast to the R implementation, the degrees-of-freedom **df** here correspond to the degree of the B-spline. The degrees-of-freedom for each effect can be specified separately in the **training\_parameters** argument of **Sddr**.

### 3.3. Functionality

Once the model is defined, the user is given several options to train, tweak, and evaluate the model. The **Sddr** class includes the functions **train**, for training the model, **eval**, for plotting the structured non-linear effects, **coeff**, for accessing the structured coefficients, and **predict**, to predict on new data. There is also the possibility to **save** and **load** the **Sddr** model. After loading a model, it is possible to resume training by adding the optional argument **resume** to

the `train` function. This will start training using the stored weights and previous optimizer state (e.g., previous learning rate and momentum).

### Case study: Airbnb modeling in Python

We first save the `airbnb` data from before as csv file.

```
R> write.csv(airbnb, file="tabular_airbnb.csv")
```

The following code demonstrates how to prepare and fit an SDDR model with **PySDDR**. We first prepare the data using the saved csv file.

```
python> import os
import numpy as np
python> mydir = os.getcwd()
python> data_path = mydir + '/data/tabular_airbnb.csv'
python> image_path = mydir + '/data/pictures/32'
python> data = pd.read_csv(data_path, delimiter=',')
python> data[['image']] = [f'{data.id[i]}.jpg' for i in range(data.shape[0])]
python> data['valid_image'] =
+ [os.path.isfile(os.path.join(image_path,
+ data.image[i]))
+ for i in range(data.shape[0])
+ ]
python> data = data[data.valid_image]
python> data = data[data['price'] > 0]
python> data.price = np.log(data.price)
```

We then define our distribution, formulas, DNN(s), training parameters, the unstructured data source, and the output directory.

```
python> distribution = 'Normal'

python> formulas = {
+ 'loc': '~ -1 + room_type + bedrooms + beds + dnn(image)',
+ 'scale': '~1 + room_type'
+ }

python> deep_models_dict = {
+ 'dnn': {
+ 'model': nn.Sequential(
+ nn.Conv2d(
+ in_channels=3,
+ out_channels=1,
+ kernel_size=(3,3),
+ stride=(4,4)
+ ),
+ nn.ReLU(),
```

```
+     nn.BatchNorm2d(num_features=1),
+     nn.MaxPool2d(kernel_size=(2,2)),
+     nn.Dropout2d(p=0.25),
+     nn.Flatten(1, -1),
+     nn.Linear(in_features=25**2, out_features=25),
+     nn.ReLU(),
+     nn.BatchNorm1d(num_features=25),
+     nn.Dropout(p=0.5)
+ ),
+ 'output_shape': 25},
+ }
```

```
python> train_parameters = {
+ 'batch_size': 32,
+ 'epochs': 100,
+ 'degrees_of_freedom': {'loc':9, 'scale':9},
+ 'val_split': 0.1,
+ 'optimizer' : optim.Adam,
+ 'early_stop_epsilon': 0.00001
+ }
```

```
python> unstructured_data = {
+ 'image' : {
+   'path' : image_path,
+   'datatype' : 'image'
+ }
+ }
```

```
python> output_dir = './outputs'
```

The SDDR model can now be initialized via:

```
python> sddr = Sddr(
+   output_dir=output_dir,
+   distribution=distribution,
+   formulas=formulas,
+   deep_models_dict=deep_models_dict,
+   train_parameters=train_parameters,
+ )
```

Finally, the model is trained using

```
python> sddr.train(
+   structured_data=data,
+   target="price",
+   unstructured_data = unstructured_data
+ )
```

Predictions can be evaluated through

```
python> data_pred = data.loc[0:1000,:]
python> ground_truth = data.loc[0:1000,'price']

python> distribution_layer, partial_effect = sddr.predict(
+ data_pred,
+ clipping=True,
+ plot=False,
+ unstructured_data = unstructured_data
+ )

python> predicted_mean = distribution_layer.loc[:,:].T
python> np.corrcoef(ground_truth, predicted_mean)
array([[1.          , 0.66732938],
       [0.66732938, 1.          ]])
```

As we can see, the prediction are correlated with the ground truth. We finally define the model `mod_simple` fitted in **deepregression** in the beginning of this tutorial and compare the results to an equivalent model definition in **PySDDR**.

```
python> tp = 'te(spline(latitude, bs="bs", df=6),' +
' spline(longitude, bs="bs", df=6))'
python> formulas = {
'loc': '~1 + ' + tp,
'scale': '~1'
}

python> train_parameters = {
'batch_size': 32,
'epochs': 1000,
'degrees_of_freedom': {'loc':9, 'scale':9},
'val_split': 0.1,
'optimizer' : optim.Adam,
'early_stop_epsilon': 0.001,
'early_stop_epochs': 10
}

python> sddr = Sddr(
output_dir=output_dir,
distribution=distribution,
formulas=formulas,
deep_models_dict=deep_models_dict,
train_parameters=train_parameters,
)

python> sddr.train(
```

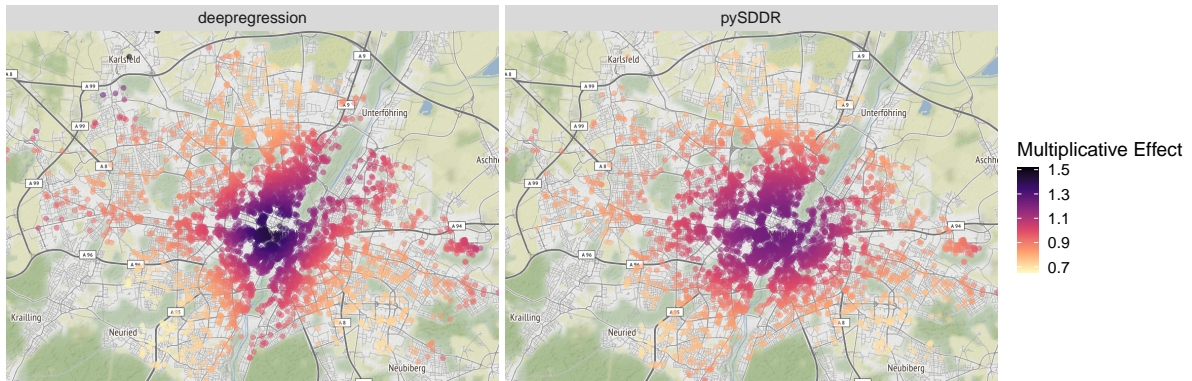


Figure 8: Comparison of the estimated effect of longitude and latitude from the two packages.

```

structured_data=data,
target="price",
unstructured_data = unstructured_data,
plot=True
)

```

We therefore extract the tensor product design matrix as follows:

```

python> data_prep = sddr.dataset[:]["datadict"]
python> desmat_tp = data_prep['loc']['structured'][:,1:].numpy()

```

and the respective coefficients using the following command:

```

python> coefs_tp = sddr.coeff('loc')[tp]

```

Figure 8 (based on codes provided in the Appendix) shows the comparison of estimated spatial effects.

For further details about **PySDDR**, we refer to the online documentation, available at <https://github.com/HelmholtzAI-Consultants-Munich/PySDDR>.

## 4. Conclusion

We have presented **deepregression**, a library that implements the SDDR framework in both R and Python. **deepregression** allows to learn the parameters of a wide range of distributions based on a combination of additive regression models and deep neural networks. The implementation is comprised of three essential components: (1) a modular neural network building system for the seamless combination of statistical and deep learning approaches, (2) an orthogonalization cell (as put forward by [Rügamer et al. \(2020\)](#)) to allow for an interpretable combination of different subnetworks, and (3) comprehensive pre-processing tools for model initialization. The software provides a common formula interface to specify distributional parameters, allows to learn shared weights and automatically handles identifiability issues. We posit that **deepregression**'s modular design and functionality provides a unique resource

for rapid and reproducible prototyping of complex statistical and deep learning models while simultaneously retaining the indispensable interpretability of classical statistical models.

## Acknowledgements

This work has been partially supported by the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036A. The authors of this work take full responsibilities for its content. Nadja Klein gratefully acknowledges support by the German research foundation (DFG) through the Emmy Noether grant KL 3037/1-1.

## References

- Baumann PFM, Hothorn T, Rügamer D (2020). “Deep Conditional Transformation Models.” *arXiv preprint arXiv:2010.07860*. [2010.07860](#).
- Blundell C, Cornebise J, Kavukcuoglu K, Wierstra D (2015). “Weight Uncertainty in Neural Networks.” *Proceedings of the 32nd International Conference on Machine Learning*, **37**, 1613–1622.
- Dunson DB (2010). “Nonparametric Bayes Applications to Biostatistics.” *Bayesian nonparametrics*, **28**, 223–273.
- Fritz C, Dorigatti E, Rügamer D (2021). “Combining Graph Neural Networks and Spatio-temporal Disease Models to Predict COVID-19 Cases in Germany.” [2101.00661](#).
- Hothorn T, Kneib T, Bühlmann P (2014). “Conditional Transformation Models.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **76**(1), 3–27. [doi: 10.1111/rssb.12017](#).
- Jacobs RA, Jordan MI, Nowlan SJ, Hinton GE (1991). “Adaptive Mixtures of Local Experts.” *Neural computation*, **3**(1), 79–87.
- Klein N, Kneib T, Lang S, Sohn A (2015). “Bayesian structured additive distributional regression with an application to regional income inequality in Germany.” *Annals of Applied Statistics*, **9**(2), 1024–1052.
- Koenker R (2005). *Quantile Regression*. Cambridge University Press, New York. Economic Society Monographs.
- Koenker R, Bassett G (1978). “Regression Quantiles.” *Econometrica*, **46**, 33–50.
- Kopper P, Pölsterl S, Wachinger C, Bischl B, Bender A, Rügamer D (2021). “Semi-Structured Deep Piecewise Exponential Models.” *AAAI Spring Symposium 2021*. [arXiv:2011.05824](#).
- Newey WK, Powell JL (1987). “Asymmetric Least Squares Estimation and Testing.” *Econometrica*, **55**, 819–847.
- Rigby RA, Stasinopoulos DM (2005). “Generalized Additive Models for Location, Scale and Shape.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **54**(3), 507–554.

- Rügamer D, Kolb C, Klein N (2020). “A Unified Network Architecture for Semi-Structured Deep Distributional Regression.” *arXiv preprint arXiv:2002.05777*. 2002.05777.
- Ruppert D, Wand MP, Carroll RJ (2003). *Semiparametric Regression*. Cambridge University Press, Cambridge and New York.
- Rügamer D, Pfisterer F, Bischl B (2020). “Neural Mixture Distributional Regression.” *arXiv preprint arXiv:2010.06889*. 2010.06889.
- Schnabel SK, Eilers P (2009). “Optimal Expectile Smoothing.” *Computational Statistics & Data Analysis*, **53**, 4168–4177.
- Tran D, Mike D, van der Wilk M, Hafner D (2019). “Bayesian Layers: A Module for Neural Network Uncertainty.” *33rd Conference on Neural Information Processing System (NeurIPS)*.
- Wood SN (2001). *mgcv: GAMs and Generalized Ridge Regression for R*.
- Wood SN (2017). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC.
- Zou H, Hastie T (2005). “Regularization and variable selection via the elastic net.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **67**(2), 301–320.

## Appendix

### Text data pre-processing

Given the uncleaned data `d` of Munich rent prices (also available at <https://github.com/davidruegamer/airbnb>), the following code generates the cleaned data used in this tutorial paper.

```
R> library(tidytext)
R> library(tm)
R> library(keras)

R> nr_words = 1000
R> maxlen = 50

R> tokenizer <- text_tokenizer(num_words = nr_words)

R> data("stop_words")
R> stopwords_regex = paste(c(stopwords('en'), stop_words$word),
  collapse = '\\b|\\b')
R> stopwords_regex = paste0('\\b', stopwords_regex, '\\b')
R> d$description = tolower(d$description)
R> d$description = stringr::str_replace_all(d$description, stopwords_regex, '')
R> d$description = gsub('[[punct:]]+', ' ', d$description)
```



```

R> stopwords_regex = paste(c(stopwords('de'), stop_words$word),
  collapse = '\\b|\\b')
R> stopwords_regex = paste0('\\b', stopwords_regex, '\\b')
R> d$description = tolower(d$description)
R> d$description = stringr::str_replace_all(d$description, stopwords_regex, "")
R> d$description = gsub('[[:punct:]]+', ' ', d$description)

R> tokenizer %>% fit_text_tokenizer(d$description)

R> text_seqs <- texts_to_sequences(tokenizer, d$description)

R> text_padded <- text_seqs %>%
  pad_sequences(maxlen = maxlen)

R> words <- data_frame(
  word = names(tokenizer$word_index),
  id = as.integer(unlist(tokenizer$word_index))
)

R> words <- words %>%
  filter(id <= tokenizer$num_words) %>%
  arrange(id)

R> saveRDS(words, file="munich_words.RDS")
R> rm(words)
R> gc()

R> text_embd <- list(texts = array(text_padded, dim=c(NROW(d), maxlen)))

R> d <- append(d, text_embd)

R> saveRDS(d, file="munich_clean_text.RDS")

```

## Map plotting

The following code creates the `map_plot` used to plot the maps of Munich.

```

R> library(ggplot2)
R> library(ggmap)
R> library(gridExtra)
R> library(viridis)
R> library(tidyr)

R> map_plot <- function(mod, d, i, legend = FALSE, python_fit = NULL){

```

```

pe_te <- get_partial_effect(mod, name = "latitude", newdata = d, which_param=i)
prices_loc <- cbind(price = pe_te[,1], d %>% select(latitude, longitude))
if(!is.null(python_fit)){

  prices_loc$price_py <- python_fit
  prices_loc <- prices_loc %>% pivot_longer(cols = c("price", "price_py"))
  prices_loc$package <- factor(prices_loc$name,
                              levels = unique(prices_loc$name),
                              labels = c("deepregression", "PySDDR"))
  prices_loc$price <- prices_loc$value

}

myLocation<-c(min(d$longitude), min(d$latitude), max(d$longitude), max(d$latitude))
myMap <- get_map(location = myLocation,
                 source="google",
                 maptype="roadmap",
                 crop=TRUE)
gg <- ggmap(myMap) + # xlab("Longitude") + ylab("Latitude") +
  geom_point(data = prices_loc, aes(x = longitude, y = latitude,
                                   colour = exp(price), alpha=0.005)) +
  scale_colour_viridis_c(option = 'magma', direction = -1,
                         name = "Multiplicative Effect") +
  guides(alpha = FALSE) + #ggtitle("Geographic Location Effect") +
  theme(plot.title = element_text(hjust = 0.5),
        text = element_text(size=10),
        axis.title.x=element_blank(),
        axis.text.x=element_blank(),
        axis.ticks.x=element_blank(),
        axis.title.y=element_blank(),
        axis.text.y=element_blank(),
        axis.ticks.y=element_blank())

if(!legend) gg <- gg + guides(colour=FALSE)
if(!is.null(python_fit)) gg <- gg + facet_grid(~ package)

return(gg)

}

}

```

## Comparison

The following create the comparison of both packages based on the codes that have been provided in the main part of this article.

```
python> partial_effect_tp = np.matmul(desmat_tp, coefs_tp)
python> np.savetxt("pe_pysddr.csv", partial_effect_tp, delimiter=",")
R> mod_simple %>% fit(
  epochs = 1000,
  early_stopping = TRUE,
  patience = 10,
  verbose = FALSE)
R> map_plot(
  mod_simple,
  as.data.frame(airbnb),
  1,
  TRUE,
  read.csv("pe_pysddr.csv", header=FALSE)$V1
)
```

**Affiliation:**

David Rügamer  
Chair of Statistical Learning & Data Science  
Department of Statistics  
LMU Munich  
80539, Munich, Germany  
E-mail: [david.ruegamer@stat.uni-muenchen.de](mailto:david.ruegamer@stat.uni-muenchen.de)