

# Preprocessing data with *recipes* : : CHEAT SHEET



## Basics

Get your data ready for modeling using 'pipable' sequences of feature engineering steps with **recipes**

```
#----- Initialize the recipe and add steps -----
rec <- recipe(x ~ ., data = train_data) |>
  step_normalize(all_numeric_predictors())

#----- Run the steps using training data -----
pr <- prep(rec, training = train_data)

#----- Apply estimates to new data -----
bake(pr, new_data = new_data)
```

**recipe**(x, ...) - Begins a new recipe specification.  
**prep**(x, ...) - Prepares the recipe with training data.  
**bake**(object, ...) - Applies estimates from prep().  
**update**(object, ...) - Updates and re-fits a model.

### COMMON `STEP\_` ARGUMENTS

**recipe** A recipe object. New steps are appended to the recipe.  
 ... Arguments passed to the external R function accessed by the step function  
**options** Selector functions to choose variables for this step

## Filters

### step\_\*(recipe, ...)

**step\_nzv**(freq\_cut = 95/5, unique\_cut = 10, options = list(freq\_cut = 95/5)) - Removes variables that are highly sparse and unbalanced.

**step\_zv**(group = NULL) - Removes variables that contain only a single value.

**step\_lincomb**(max\_steps = 5) - Removes numeric variables that have exact linear combinations between them.

**step\_corr**(threshold = 0.9, use = "pairwise.complete.obs", method = "pearson") - Removes variables that have large absolute correlations with other variables.

**step\_filter\_missing**(threshold = 0.1) - Removes variables that have too many missing values.

**step\_rm**() - Removes selected variables.

## In-place Transformations

**step\_mutate**(recipe, ..., .pkgs = character()) - General purpose transformer using dplyr.  
**step\_relu**(recipe, ..., shift = 0, reverse = FALSE, smooth = FALSE, prefix = "right\_relu\_") - Applies smoothed rectified linear transformation.  
**step\_sqrt**(recipe, ...) - Applies square root transformation.

### BASIS FUNCTIONS

#### step\_\*(recipe, ..., keep\_original\_cols = FALSE)

**step\_spline\_natural**(deg\_free = 10, options = NULL) - Creates a natural spline (a.k.a restricted cubic spline) features.

**step\_spline\_b**(deg\_free = 10, degree = 3, options = NULL) - Creates b-spline features.  
**step\_spline\_convex**(deg\_free = 10, degree = 3, options = NULL)

**step\_spline\_monotone**(deg\_free = 10, degree = 3, options = NULL)

**step\_spline\_nonnegative**(deg\_free = 10, degree = 3, options = NULL)

**step\_poly**(degree = 2L, options = list()) - Creates new columns that are basis expansions of variables using orthogonal polynomials.

**step\_poly\_bernstein**(degree = 10, options = NULL, results = NULL) - Creates Bernstein polynomial features.

### NORMALIZATION

#### step\_\*(recipe, ...)

**step\_normalize**(na\_rm = TRUE) - Normalizes to have a standard deviation of 1 and mean of 0.

**step\_YeoJohnson**() - Makes data look more like a normal distribution.

**step\_percentile**(options = list(probs = (0:100)/100), outside = "none") - Replaces the value of a variable with its percentile from the training set.

**step\_range**(min = 0, max = 1, clipping = TRUE) - Normalizes numeric data to be within a pre-defined range of values.

**step\_spatialsign**(na\_rm = TRUE) - Converts numeric data into a projection on to a unit sphere.

### DISCRETIZE

#### step\_\*(recipe, ...)

**step\_discretize**(num\_breaks = 4, min\_unique = 10, options = list(prefix = "bin")) - Converts numeric data into a factor with bins having approximately the same number of data points.

**step\_cut**(breaks, include\_outside\_range = FALSE) - Cuts a numeric variable into a factor based on provided boundary values.

## Imputation

#### step\_\*(recipe, ...)

**step\_impute\_bag**(impute\_with = all\_predictors(), trees = 25, options = list(keepX = FALSE)) - Creates a bagged tree model for data. Good for categorical data.

**step\_impute\_knn**(neighbors = 5, impute\_with = all\_predictors(), options = list(nthread = 1, eps = 1e-08)) - Uses Gower's distance which can be used for mixtures of nominal and numeric data.

**step\_impute\_linear**(impute\_with = all\_predictors()) - Creates linear regression models to impute missing data.

**step\_impute\_lower**(threshold = NULL) - Substitutes the truncated value by a random number between zero and the truncation point.

**step\_impute\_mean**(trim = 0) - Substitutes missing values of numeric variables by the training set mean of those variables.

**step\_impute\_median**() - Substitutes missing values of numeric variables by the training set median of those variables.

**step\_impute\_mode**() - Imputes nominal data using the most common value.

**step\_impute\_roll**(statistic = median, window = 5L) - Imputes numeric data using a rolling window statistic.

**step\_unknown**(new\_level = "unknown") - Assigns a missing value in a factor level to "unknown".

## Encodings

### step\_\*(recipe, ...)

#### TYPE CONVERTERS



**step\_factor2string**() - Converts one or more factor vectors to strings.

**step\_string2factor**() - Converts one or more character vectors to factors (ordered or unordered).

**step\_num2factor**(transform = function(x) x) - Converts one or more numeric vectors to factors (ordered or unordered). This can be useful when categories are encoded as integers.



**step\_integer**(strict = TRUE, zero\_based = FALSE) - Converts data into a set of ascending integers based on the ascending order from the training data.



#### VALUE CONVERTERS



**step\_indicate\_na**(sparse = "auto", keep\_original\_cols = TRUE) - Creates and append additional binary columns to the data set to indicate which observations are missing.



**step\_ordinalscore**(convert = as.numeric) - Converts ordinal factor variables into numeric scores.



**step\_unorder**() - Turns ordered factor variables into unordered factor variables.

#### OTHER



**step\_relevel**(ref\_level) - Reorders factor columns so that the level specified by ref\_level is first. This is useful for contr.treatment() contrasts which take the first level as the reference.



**step\_novel**(new\_level = "new") - Assigns a previously unseen factor level to "new".



**step\_other**(threshold = 0.05, other = "other") - Pools infrequently occurring values into an "other" category.

