

---

# DataWig Documentation

**Amazon**

**Jul 31, 2020**



<b>1</b>	<b>Table of Contents</b>	<b>3</b>
1.1	Introduction . . . . .	3
1.2	User Guide . . . . .	3
1.2.1	Step-by-step Examples . . . . .	3
1.2.2	Introduction to SimpleImputer . . . . .	5
1.2.3	Introduction to Imputer . . . . .	7
1.2.4	Parameters for Different Data Types . . . . .	8
1.2.5	Advanced Features . . . . .	9
1.3	Contributing to DataWig . . . . .	10
1.4	API . . . . .	10
1.4.1	Simple Imputer . . . . .	10
1.4.2	Imputer . . . . .	15
1.4.3	Column Encoders . . . . .	18
	<b>Python Module Index</b>	<b>25</b>
	<b>Index</b>	<b>27</b>



This is the documentation for DataWig, a framework for learning models to impute missing values in tables.



## 1.1 Introduction

This is the documentation for **DataWig**, a framework for learning models to impute missing values in tables.

Details on the underlying model can be found in [Biessmann, Salinas et al. 2018](#).

```
@inproceedings{datawig,  
  author = {Biessmann, Felix and Salinas, David and Schelter, Sebastian and Schmidt,   
↪Philipp and Lange, Dustin},  
  title = {"Deep" Learning for Missing Value Imputationin Tables with Non-Numerical   
↪Data},  
  booktitle = {Proceedings of the 27th ACM International Conference on Information and   
↪Knowledge Management},  
  series = {CIKM '18},  
  year = {2018},  
  isbn = {978-1-4503-6014-2},  
  location = {Torino, Italy},  
  pages = {2017--2025},  
  numpages = {9},  
  url = {http://doi.acm.org/10.1145/3269206.3272005},  
  doi = {10.1145/3269206.3272005},  
  keywords = {data cleaning, missing value imputation}}
```

## 1.2 User Guide

### 1.2.1 Step-by-step Examples

#### Setup

For installing DataWig, follow the installation instructions in the [readme](#).

### Examples

In each example, we provide a detailed description of important features along with python code that highlights these features on a public dataset. We recommend reading through the overview of DataWig and then following the below examples in order.

For additional examples and use cases, refer to the [unit test cases](#).

### Data

Unless otherwise specified, these examples will make use of the [Multimodal Attribute Extraction \(MAE\) dataset](#). This dataset contains over 2.2 million products with corresponding attributes, but to make data loading and processing more manageable, we provide a reformatted subset of the validation data (for the *finish* and *color* attributes) as a .csv file.

This data contains columns for *title*, *text*, *finish*, and *color*. The title and text columns contain string data that will be used to impute the finish attribute. Note, the dataset is extremely noisy, but still provides a good example for real-world use cases of DataWig.

To speed up run-time, all examples will use a smaller version of this finish dataset that contains ~5000 samples. Run the following in this directory to download this dataset:

```
wget https://github.com/aws-labs/datawig/raw/master/examples/mae_train_dataset.csv
```

To get the complete finish dataset with all data, please check instructions [here](#).

If you'd like to use this data in your own experiments, please remember to cite the original MAE paper:

```
@article{RobertLLogan2017MultimodalAE,  
  title={Multimodal Attribute Extraction},  
  author={IV RobertL.Logan and Samuel Humeau and Sameer Singh},  
  journal={CoRR},  
  year={2017},  
  volume={abs/1711.11118}}
```

### Overview of DataWig

Here, we give a brief overview of the internals of DataWig.

#### ColumnEncoder (*column\_encoder.py*)

Defines an abstract super class of column encoders that transforms the raw data of a column (e.g. strings from a product title) into an encoded numerical representation.

There are a few options for ColumnEncoders (subclasses) depending on the column data type:

- `SequentialEncoder`: for sequences of string symbols (e.g. characters or words)
- `BowEncoder`: bag-of-word representation for strings, as sparse vectors
- `CategoricalEncoder`: for categorical variables (one-hot encoding)
- `NumericalEncoder`: for numerical values



## Featurizer (*mxnet\_input\_symbol.py*)

Defines a specific featurizer for data that has been encoded into a numerical format by `ColumnEncoder`. The Featurizer is used to feed data into the imputation model's computational graph for training and prediction.

There are a few options for Featurizers depending on which `ColumnEncoder` was used for a particular column:

- `LSTMFeaturizer` maps an input representing a sequence of symbols into a latent vector using an LSTM
- `BowFeaturizer` used with `BowEncoder` on string data
- `EmbeddingFeaturizer` maps encoded categorical data into a vector representations (word-embeddings)
- `NumericalFeaturizer` extracts features from numerical data using fully connected layers

## SimpleImputer (*simple\_imputer.py*)

Using `SimpleImputer` is the easiest way to deploy an imputation model on your dataset with DataWig. As the name suggests, the `SimpleImputer` is straightforward to call from a python script and uses default encoders and featurizers that usually yield good results on a variety of datasets.

## Imputer (*imputer.py*)

`Imputer` is the backbone of the `SimpleImputer` and is responsible for running the preprocessing code, creating the model, executing training, and making predictions. Using the `Imputer` enables more flexibility with specifying model parameters, such as using particular encoders and featurizers rather than the default ones that `SimpleImputer` uses.

## 1.2.2 Introduction to SimpleImputer

This tutorial will teach you the basics of how to use `SimpleImputer` for your data imputation tasks. As an advanced feature the `SimpleImputer` supports label-shift detection and correction which is described in *Label Shift and Empirical Risk Minimization*. For now, we will use a subset of the MAE data as an example. To download this data, please refer to the previous section.

Open the [SimpleImputer intro](#) in this directory to see the code used in this tutorial.

### Load Data

First, let's load the data into a pandas `DataFrame` and split the data into train (80%) and test (20%) subsets.

```
df = pd.read_csv('../finish_val_data_sample.csv')
df_train, df_test = random_split(df, split_ratios=[0.8, 0.2])
```

Note, the `random_split()` method is provided in `datawig.utils`. The validation set is partitioned from the train data during training and defaults to 10%.

### Default SimpleImputer

At the most basic level, you can run the `SimpleImputer` on data without specifying any additional arguments. This will automatically choose the right `ColumnEncoder` and `Featurizer` for each column and train an imputation model with default hyperparameters.

To train a model, you can simply initialize a `SimpleImputer`, specifying the input columns containing useful data for imputation, the output column that you'd like to impute values for, and the output path, which will store model data and metrics. Then, you can use the `fit()` method to train the model.

```
#Initialize a SimpleImputer model
imputer = SimpleImputer(
    input_columns=['title', 'text'],
    output_column='finish',
    output_path = 'imputer_model'
)

#Fit an imputer model on the train data
imputer.fit(train_df=df_train)
```

From here, you can use this model to make predictions on the test set and return the original dataframe with an additional column containing the model's predictions.

```
predictions = imputer.predict(df_test)
```

Finally, you can determine useful metrics to gauge how well the model's predictions compare to the true values (using `sklearn.metrics`).

```
#Calculate f1 score
f1 = f1_score(predictions['finish'], predictions['finish_imputed'])

#Print overall classification report
print(classification_report(predictions['finish'], predictions['finish_imputed']))
```

### HPO with SimpleImputer

DataWig also enables hyperparameter optimization to find the best model on a particular dataset.

The steps for training a model with HPO are identical to the default `SimpleImputer`.

```
imputer = SimpleImputer(
    input_columns=['title', 'text'],
    output_column='finish',
    output_path='imputer_model'
)

# fit an imputer model with customized hyperparameters
imputer.fit_hpo(train_df=df_train)
```

Calling HPO like this will search through some basic and usually helpful hyperparameter choices. There are two ways for a more detailed search. Firstly, `fit_hpo` offers additional arguments that can be inspected in the `SimpleImputer`. For even more configurations and variation of hyperparameters for the various input column types, a dictionary with ranges can be passed to `fit_hpo` as can be seen in the [hpo-code](#). Results for any HPO run can be accessed under `imputer.hpo.results` and the model from any HPO run can then be loaded using `imputer.load_hpo_model(idx)` passing the model index.

### Load Saved Model

Once a model is trained, it will be saved in the location of `output_path`, which you specified as an argument when initializing the `SimpleImputer`. You can easily load this model for further experiments or run on new datasets as follows.

```
#Load saved model
imputer = SimpleImputer.load('./imputer_model')
```

This model also contains the associated metrics (stored as a dictionary) calculated on the validation set during training.

```
#Load metrics from the validation set
metrics = imputer.load_metrics()
weighted_f1 = metrics['weighted_f1']
avg_precision = metrics['avg_precision']
# ...
```

### 1.2.3 Introduction to Imputer

This tutorial will teach you the basics of how to use the `Imputer` for your data imputation tasks. We will use a subset of the MAE data as an example. To download this data, please refer to [README](#).

Open [Imputer intro](#) to see the code used in this tutorial.

#### Load Data

First, let's load the data into a pandas `DataFrame` and split the data into train (80%) and test (20%) subsets.

```
df = pd.read_csv('../finish_val_data_sample.csv')
df_train, df_test = random_split(df, split_ratios=[0.8, 0.2])
```

Note, the `random_split()` method is provided in `datawig.utils`. The validation set is partitioned from the train data during training and defaults to 10%.

#### Default Imputer

The key difference with the `Imputer` is specifying the Encoders and Featurizers used for particular columns in your dataset. Once this is done, initializing the model, training, and making predictions with the `Imputer` is similar to the `SimpleImputer`

```
#Specify encoders and featurizers
data_encoder_cols = [BowEncoder('title'), BowEncoder('text')]
label_encoder_cols = [CategoricalEncoder('finish')]
data_featurizer_cols = [BowFeaturizer('title'), BowFeaturizer('text')]

imputer = Imputer(
    data_featurizers=data_featurizer_cols,
    label_encoders=label_encoder_cols,
    data_encoders=data_encoder_cols,
    output_path='imputer_model'
)

imputer.fit(train_df=df_train)
predictions = imputer.predict(df_test)
```

For the input columns that contain data useful for imputation, the `Imputer` expects you to specify the particular encoders and featurizers. For the label column that you are trying to impute, only specifying the type of encoder is necessary.

### Using Different Encoders and Featurizers

One of the key advantages with the `Imputer` is that you get flexibility for customizing exactly which encoders and featurizers to use, which is something you can't do with the `SimpleImputer`.

For example, let's say you wanted to use an LSTM rather than the default bag-of-words text model that the `SimpleImputer` uses. To do this, you can simply specify the proper encoders and featurizers to initialize the `Imputer` model.

```
#Using LSTMs instead of bag-of-words
data_encoder_cols = [SequentialEncoder('title'), SequentialEncoder('text')]
label_encoder_cols = [CategoricalEncoder('finish')]
data_featurizer_cols = [LSTMFeaturizer('title'), LSTMFeaturizer('text')]

imputer = Imputer(
    data_featurizers=data_featurizer_cols,
    label_encoders=label_encoder_cols,
    data_encoders=data_encoder_cols,
    output_path='imputer_model'
)
```

### Prediction with Probabilities

Beyond directly predicting values, the `Imputer` can also return the probabilities for each class on every sample (numpy array of shape samples-by-labels). This can help with understanding what the model is predicting and with what probability for each sample.

```
prob_dict = imputer.predict_proba(df_test)
```

In addition, you can get the probabilities only for the top-k most likely predicted classes (rather than for all the classes above).

```
prob_dict_topk = imputer.predict_proba_top_k(df_test, top_k=5)
```

### Get Predictions and Metrics

To get predictions (original dataframe with an extra column) and the associated metrics from the validation set during training, you can run the following:

```
predictions, metrics = imputer.transform_and_compute_metrics(df_test)
```

## 1.2.4 Parameters for Different Data Types

This tutorial will highlight the different parameters associated with column data types supported by DataWig. We use the `SimpleImputer` in these examples, but the same concepts apply when using the `Imputer` and other encoders/featurizers.

The [parameter tutorial](#) contains the complete code for training models on text and numerical data. Here, we illustrate examples of relevant parameters for training models on each of these types of data.

It's important to note that your dataset can contain columns with mixed types. The `SimpleImputer` automatically determines which encoder and featurizer to use when training an imputation model!

## Text Data

The key parameters associated with text data are:

- `num_hash_buckets` dimensionality of the vector for bag-of-words
- `tokens` type of tokenization used for text data (default: chars)

Here is an example of using these parameters:

```
imputer_text.fit_hpo(
    train_df=df_train,
    num_epochs=50,
    learning_rate_candidates=[1e-3, 1e-4],
    final_fc_hidden_units_candidates=[[100]],
    num_hash_bucket_candidates=[2**10, 2**15],
    tokens_candidates=['chars', 'words']
)
```

Apart from the text parameters, `final_fc_hidden_units` corresponds to a list containing the dimensionality of the fully connected layer after all column features are concatenated. The length of this list is the number of hidden fully connected layers.

## Numerical Data

The key parameters associated with numerical data are:

- `latent_dim` dimensionality of the fully connected layers for creating a feature vector from numerical data
- `hidden_layers` number of fully connected layers

Here is an example of using these parameters:

```
imputer_numeric.fit_hpo(
    train_df=df_train,
    num_epochs=50,
    learning_rate_candidates=[1e-3, 1e-4],
    latent_dim_candidates=[50, 100],
    hidden_layers_candidates=[0, 2],
    final_fc_hidden_units=[[100]]
)
```

In this case, the model will use a fully connected layer size of 50 or 100, with 0 or 2 hidden layers.

## 1.2.5 Advanced Features

### Label Shift and Empirical Risk Minimization

The `SimpleImputer` implements the method described by [Lipton, Wang and Smola](#) to detect and fix label shift for categorical outputs. Label shift occurs when the marginal distribution differs between the training and production setting. For instance, we might be interested in imputing the color of T-Shirts from their free-text description. Let's assume that the training data consists only of women's T-Shirts while the production data consists only of Men's T-Shirts. Then the marginal distribution of colors,  $p(\text{color})$ , is likely different while the conditional,  $p(\text{description} | \text{color})$  may be unchanged. This is a scenario where `datawig` can detect and fix the shift.

Upon training a `SimpleImputer`, we can detect shift by calling:

```
weights = imputer.check_for_label_shift(production_data)
```

Note, that `production_data` needs to have all the relevant input columns but does not have labels. This call will log the severity of the shift and further information, as follows.

```
The estimated true label marginals are [('black', 0.62), ('white', 0.38)]
Marginals in the training data are [('black', 0.23), ('white', 0.77)]
Reweighting factors for empirical risk minimization{'label_0': 2.72, 'label_1': 0.49}
The smallest eigenvalue of the confusion matrix is 0.21 ' (needs to be > 0).
```

To fix the shift, the reweighting factors are most important. They are returned as dictionary where each key is a label and the value is the corresponding weight by which any observation's contribution to the log-likelihood must be multiplied to minimize the empirical risk. To correct the shift we need to retrain the model with a weighted likelihood which can easily be achieved by passing the weight dictionary to the `fit()` method.

```
simple_imputer.fit(train_df, class_weights=weights)
```

The resulting model will generally have improved performance on the `production_data`, if there was a label shift present and if the original classifier performed reasonably well. For further assumptions see the above cited paper. Note, that in extreme cases such as very high label noise, this method can lead to a decreased model performance.

Reweighting the likelihood can be useful for reasons other than label-shift. For instance we may trust certain observations more than others and wish to up-weight their impact on the model parameters. To this end, weights can also be passed on an instance level as list with an entry for every row in the training data, for instance:

```
simple_imputer.fit(train_df, class_weights=[1, 1, 1, 2, 1, 1, 1, ...])
```

## 1.3 Contributing to DataWig

Please see [how to contribute](#).

## 1.4 API

### 1.4.1 Simple Imputer

DataWig SimpleImputer: Uses some simple default encoders and featurizers that usually yield decent imputation quality

```
class datawig.simple_imputer.SimpleImputer (input_columns: List[str], output_column: str,  
output_path: str = "", num_hash_buckets: int  
= 32768, num_labels: int = 100, tokens: str  
= 'chars', numeric_latent_dim: int = 100, nu-  
meric_hidden_layers: int = 1, is_explainable:  
bool = False)
```

SimpleImputer model based on n-grams of concatenated strings of input columns and concatenated numerical features, if provided.

Given a data frame with string columns, a model is trained to predict observed values in label column using values observed in other columns.

The model can then be used to impute missing values.

#### Parameters

- **input\_columns** – list of input column names (as strings)
- **output\_column** – output column name (as string)
- **output\_path** – path to store model and metrics
- **num\_hash\_buckets** – number of hash buckets used for the n-gram hashing vectorizer, only used for non-numerical input columns, ignored otherwise
- **num\_labels** – number of imputable values considered after, only used for non-numerical input columns, ignored otherwise
- **tokens** – string, ‘chars’ or ‘words’ (default ‘chars’), determines tokenization strategy for n-grams, only used for non-numerical input columns, ignored otherwise
- **numeric\_latent\_dim** – int, number of latent dimensions for hidden layer of NumericalFeaturizers; only used for numerical input columns, ignored otherwise
- **numeric\_hidden\_layers** – number of numeric hidden layers
- **is\_explainable** – if this is True, a stateful tf-idf encoder is used that allows explaining classes and single instances

Example usage:

```
from datawig.simple_imputer import SimpleImputer import pandas as pd

fn_train = os.path.join(datawig_test_path, "resources", "shoes", "train.csv.gz") fn_test =
os.path.join(datawig_test_path, "resources", "shoes", "test.csv.gz")

df_train = pd.read_csv(training_data_files) df_test = pd.read_csv(testing_data_files)

output_path = "imputer_model"

# set up imputer model imputer = SimpleImputer( input_columns=['item_name', 'bullet_point'], out-
put_column='brand')

# train the imputer model imputer = imputer.fit(df_train)

# obtain imputations imputations = imputer.predict(df_test)

check_data_types (data_frame: pandas.core.frame.DataFrame) → None
    Checks whether a column contains string or numeric data
```

**Parameters** *data\_frame* –

**Returns**

```
check_for_label_shift (target_data: pandas.core.frame.DataFrame) → dict
    Detect label shift in the validation data
```

**Parameters**

- **test\_data** – data frame that contains labels
- **target\_data** – unlabelled data for which predictions are to be generated

**Returns** dictionary with labels as keys and weights as values.

```
static complete (data_frame: pandas.core.frame.DataFrame, precision_threshold: float = 0.0, in-
place: bool = False, hpo: bool = False, verbose: int = 0, num_epochs: int = 100,
iterations: int = 1, output_path: str = '.')
```

Given a dataframe with missing values, this function detects all imputable columns, trains an imputation model on all other columns and imputes values for each missing value. Several imputation iterators can be run. Imputable columns are either numeric columns or non-numeric categorical columns; for determining whether a

column is categorical (as opposed to a plain text column) we use the following heuristic: a non-numeric categorical column should have least 10 times as many rows as there were unique values

**If an imputation model did not reach the precision specified in the `precision_threshold` parameter for a given imputation value, that value will not be imputed; thus depending on the `precision_threshold`, the returned dataframe can still contain some missing values.**

For numeric columns, we do not filter for accuracy. :param data\_frame: original dataframe :param precision\_threshold: precision threshold for categorical imputations (default: 0.0) :param inplace: whether or not to perform imputations inplace (default: False) :param hpo: whether or not to perform hyperparameter optimization (default: False) :param verbose: verbosity level, values > 0 log to stdout (default: 0) :param num\_epochs: number of epochs for each imputation model training (default: 100) :param iterations: number of iterations for iterative imputation (default: 1) :param output\_path: path to store model and metrics :return: dataframe with imputations

**explain** (*label: str, k: int = 10, label\_column: str = None*) → dict

Return dictionary with a list of tuples for each explainable input column. Each tuple denotes one of the top k features with highest correlation to the label.

#### Parameters

- **label** – label value to explain
- **k** – number of explanations for each input encoder to return. If not given, return top 10 explanations.
- **label\_column** – name of label column to be explained (optional, defaults to the first available column.)

**explain\_instance** (*instance: pandas.core.series.Series, k: int = 10, label\_column: str = None, label: str = None*) → dict

Return dictionary with list of tuples for each explainable input column of the given instance. Each entry shows the most highly correlated features to the given label (or the top predicted label of not provided).

#### Parameters

- **instance** – row of data frame (or dictionary)
- **k** – number of explanations (ngrams) for text inputs
- **label\_column** – name of label column to be explained (optional)
- **label** – explain why instance is classified as label, otherwise explain top-label per input

**fit** (*train\_df: pandas.core.frame.DataFrame, test\_df: pandas.core.frame.DataFrame = None, ctx: <module 'mxnet.context' from '/home/docs/checkouts/readthedocs.org/user\_builds/datawig/envs/latest/lib/python3.7/site-packages/mxnet/context.py'> = [cpu(0)], learning\_rate: float = 0.004, num\_epochs: int = 100, patience: int = 5, test\_split: float = 0.1, weight\_decay: float = 0.0, batch\_size: int = 16, final\_fc\_hidden\_units: List[int] = None, calibrate: bool = True, class\_weights: dict = None, instance\_weights: list = None*) → Any

Trains and stores imputer model

#### Parameters

- **train\_df** – training data as dataframe
- **test\_df** – test data as dataframe; if not provided, a ratio of test\_split of the training data are used as test data
- **ctx** – List of mxnet contexts (if no gpu's available, defaults to [mx.cpu()]) User can also pass in a list gpus to be used, ex. [mx.gpu(0), mx.gpu(2), mx.gpu(4)]
- **learning\_rate** – learning rate for stochastic gradient descent (default 4e-4)



- **num\_epochs** – maximal number of training epochs (default 10)
- **patience** – used for early stopping; after [patience] epochs with no improvement, training is stopped. (default 3)
- **test\_split** – if no test\_df is provided this is the ratio of test data to be held separate for determining model convergence
- **weight\_decay** – regularizer (default 0)

:param batch\_size (default 16) :param final\_fc\_hidden\_units: list dimensions for FC layers after the final concatenation :param calibrate: Control automatic model calibration :param class\_weights: Dictionary with labels as keys and weights as values.

Weighs each instance's contribution to the likelihood based on the corresponding class.

**Parameters instance\_weights** – List of weights for each instance in train\_df.

```
fit_hpo (train_df: pandas.core.frame.DataFrame, test_df: pandas.core.frame.DataFrame =
None, hps: dict = None, num_evals: int = 10, max_running_hours: float = 96.0,
hpo_run_name: str = None, user_defined_scores: list = None, num_epochs: int =
None, patience: int = None, test_split: float = 0.2, weight_decay: List[float] = None,
batch_size: int = 16, num_hash_bucket_candidates: List[float] = [4096, 32768, 262144],
tokens_candidates: List[str] = ['words', 'chars'], numeric_latent_dim_candidates: List[int]
= None, numeric_hidden_layers_candidates: List[int] = None, final_fc_hidden_units:
List[List[int]] = None, learning_rate_candidates: List[float] = None, normalize_numeric:
bool = True, hpo_max_train_samples: int = None, ctx: <module 'mxnet.context' from
'/home/docs/checkouts/readthedocs.org/user_builds/datawig/envs/latest/lib/python3.7/site-
packages/mxnet/context.py'> = [cpu(0)]) → Any
```

Fits an imputer model with hyperparameter optimization. The parameter ranges are searched randomly.

Grids are specified using the \*\_candidates arguments (old) or with more flexibility via the dictionary hps.

#### Parameters

- **train\_df** – training data as dataframe
- **test\_df** – test data as dataframe; if not provided, a ratio of test\_split of the training data are used as test data
- **hps** – nested dictionary where hps[global][parameter\_name] is list of parameters. Similarly, hps[column\_name][parameter\_name] is a list of parameter values for each input column. Further, hps[column\_name]['type'] is in ['numeric', 'categorical', 'string'] and is inferred if not provided.
- **num\_evals** – number of evaluations for random search
- **max\_running\_hours** – Time before the hpo run is terminated in hours.
- **hpo\_run\_name** – string to identify the current hpo run.
- **user\_defined\_scores** – list with entries (Callable, str), where callable is a function accepting \*\*kwargs true, predicted, confidence. Allows custom scoring functions.

Below are parameters of the old implementation, kept to ascertain backwards compatibility. :param num\_epochs: maximal number of training epochs (default 10) :param patience: used for early stopping; after [patience] epochs with no improvement,

training is stopped. (default 3)

#### Parameters

- **test\_split** – if no test\_df is provided this is the ratio of test data to be held separate for determining model convergence
- **weight\_decay** – regularizer (default 0)

:param batch\_size (default 16) :param num\_hash\_bucket\_candidates: candidates for gridsearch hyperparameter

optimization (default [2\*\*10, 2\*\*13, 2\*\*15, 2\*\*18, 2\*\*20])

### Parameters

- **tokens\_candidates** – candidates for tokenization (default ['words', 'chars'])
- **numeric\_latent\_dim\_candidates** – candidates for latent dimensionality of numerical features (default [10, 50, 100])
- **numeric\_hidden\_layers\_candidates** – candidates for number of hidden layers of
- **final\_fc\_hidden\_units** – list of lists w/ dimensions for FC layers after the final concatenation (NOTE: for HPO, this expects a list of lists)
- **learning\_rate\_candidates** – learning rate for stochastic gradient descent (default 4e-4) numerical features (default [0, 1, 2])
- **learning\_rate\_candidates** – candidates for learning rate (default [1e-1, 1e-2, 1e-3])
- **normalize\_numeric** – boolean indicating whether or not to normalize numeric values
- **hpo\_max\_train\_samples** – training set size for hyperparameter optimization. use is deprecated.
- **ctx** – List of mxnet contexts (if no gpu's available, defaults to [mx.cpu()]) User can also pass in a list gpus to be used, ex. [mx.gpu(0), mx.gpu(2), mx.gpu(4)] This parameter is deprecated.

**Returns** pd.DataFrame with with hyper-parameter configurations and results

**static load** (*output\_path: str*) → Any  
Loads model from output path

**Parameters** **output\_path** – output\_path field of trained SimpleImputer model

**Returns** SimpleImputer model

**load\_hpo\_model** (*hpo\_name: int = None*)  
Load model after hyperparameter optimisation has ran. Overwrites local artifacts of self.imputer.

**Parameters** **hpo\_name** – Index of the model to be loaded. Default, load model with highest weighted precision or mean squared error.

**Returns** imputer object

**load\_metrics** () → Dict[str, Any]  
Loads various metrics of the internal imputer model, returned as dictionary :return: Dict[str,Any]

**predict** (*data\_frame: pandas.core.frame.DataFrame, precision\_threshold: float = 0.0, imputation\_suffix: str = '\_imputed', score\_suffix: str = '\_imputed\_proba', inplace: bool = False*)

**Imputes most likely value if it is above a certain precision threshold determined on the** validation set

Precision is calculated as part of the `datawig.evaluate_and_persist_metrics` function.

Returns original dataframe with imputations and respective likelihoods as estimated by imputation model; in additional columns; names of imputation columns are that of the label suffixed with `imputation_suffix`, names of respective likelihood columns are suffixed with `score_suffix`

#### Parameters

- **data\_frame** – data frame (pandas)
- **precision\_threshold** – double between 0 and 1 indicating precision threshold
- **imputation\_suffix** – suffix for imputation columns
- **score\_suffix** – suffix for imputation score columns
- **inplace** – add column with imputed values and column with confidence scores to data\_frame, returns the modified object (True). Create copy of data\_frame with additional columns, leave input unmodified (False).

**Returns** data\_frame original dataframe with imputations and likelihood in additional column

#### `save()`

Saves model to disk; mxnet module and imputer are stored separately

## 1.4.2 Imputer

DataWig Imputer: Imputes missing values in tables

```
class datawig.imputer.Imputer (data_encoders: List[datawig.column_encoders.ColumnEncoder],
                               data_featurizers: List[datawig.mxnet_input_symbols.Featurizer],
                               label_encoders: List[datawig.column_encoders.ColumnEncoder],
                               output_path="")
```

Imputer model based on deep learning trained with MxNet

Given a data frame with string columns, a model is trained to predict observed values in one or more column using values observed in other columns. The model can then be used to impute missing values.

#### Parameters

- **data\_encoders** – list of `datawig.mxnet_input_symbol.ColumnEncoders`, output\_column name must match field\_name of data\_featurizers
- **data\_featurizers** – list of `Featurizer`;
- **label\_encoders** – list of `CategoricalEncoder` or `NumericalEncoder`
- **output\_path** – path to store model and metrics

```
calibrate (test_iter: datawig.iterators.ImputerIterDf)
```

Checks model calibration and fits temperature scaling. If the fit improves model calibration, the temperature parameter is assigned as property to self and used for all further predictions in `self.predict_mxnet_iter()`. Saves calibration information to dictionary.

**Parameters** `test_iter` – iterator, see `ImputerIter` in `iterators.py`

**Returns** None

```
explain (label: str, k: int = 10, label_column: str = None) → dict
```

Return dictionary with a list of tuples for each explainable input column. Each tuple denotes one of the top k features with highest correlation to the label.

#### Parameters

- **label** – label value to explain
- **k** – number of explanations for each input encoder to return. If not given, return top 10 explanations.
- **label\_column** – name of label column to be explained (optional, defaults to the first available column.)

**explain\_instance** (*instance: pandas.core.series.Series, k: int = 10, label\_column: str = None, label: str = None*) → dict

Return dictionary with list of tuples for each explainable input column of the given instance. Each entry shows the most highly correlated features to the given label (or the top predicted label of not provided).

#### Parameters

- **instance** – row of data frame (or dictionary)
- **k** – number of explanations (ngrams) for text inputs
- **label\_column** – name of label column to be explained (optional)
- **label** – explain why instance is classified as label, otherwise explain top-label per input

**fit** (*train\_df: pandas.core.frame.DataFrame, test\_df: pandas.core.frame.DataFrame = None, ctx: <module 'mxnet.context' from '/home/docs/checkouts/readthedocs.org/user\_builds/datawig/envs/latest/lib/python3.7/site-packages/mxnet/context.py'> = [cpu(0)], learning\_rate: float = 0.001, num\_epochs: int = 100, patience: int = 3, test\_split: float = 0.1, weight\_decay: float = 0.0, batch\_size: int = 16, final\_fc\_hidden\_units: List[int] = None, calibrate: bool = True*)

Trains and stores imputer model

#### Parameters

- **train\_df** – training data as dataframe
- **test\_df** – test data as dataframe; if not provided, [test\_split] % of the training data are used as test data
- **ctx** – List of mxnet contexts (if no gpu's available, defaults to [mx.cpu()]) User can also pass in a list gpus to be used, ex. [mx.gpu(0), mx.gpu(2), mx.gpu(4)]
- **learning\_rate** – learning rate for stochastic gradient descent (default 1e-4)
- **num\_epochs** – maximal number of training epochs (default 100)
- **patience** – used for early stopping; after [patience] epochs with no improvement, training is stopped. (default 3)
- **test\_split** – if no test\_df is provided this is the ratio of test data to be held separate for determining model convergence
- **weight\_decay** – regularizer (default 0)
- **batch\_size** – default 16
- **final\_fc\_hidden\_units** – list of dimensions for the final fully connected layer.
- **calibrate** – whether to calibrate predictions

**Returns** trained imputer model

**static load** (*output\_path: str*) → Any

Loads model from output path

**Parameters** **output\_path** – output\_path field of trained Imputer model

**Returns** imputer model

**predict** (*data\_frame*: *pandas.core.frame.DataFrame*, *precision\_threshold*: *float = 0.0*, *imputation\_suffix*: *str = '\_imputed'*, *score\_suffix*: *str = '\_imputed\_proba'*, *inplace*: *bool = False*)  
→ *pandas.core.frame.DataFrame*

Computes imputations for numerical or categorical values

For categorical imputations, most likely values are imputed if values are above a certain precision threshold computed on the validation set Precision is calculated as part of the *datawig.evaluate\_and\_persist\_metrics* function.

For numerical imputations, no thresholding is applied.

Returns original dataframe with imputations and respective likelihoods as estimated by imputation model in additional columns; names of imputation columns are that of the label suffixed with *imputation\_suffix*, names of respective likelihood columns are suffixed with *score\_suffix*

#### Parameters

- **data\_frame** – pandas data\_frame
- **precision\_threshold** – double between 0 and 1 indicating precision threshold for each imputation
- **imputation\_suffix** – suffix for imputation columns
- **score\_suffix** – suffix for imputation score columns
- **inplace** – add column with imputed values and column with confidence scores to data\_frame, returns the modified object (True). Create copy of data\_frame with additional columns, leave input unmodified (False).

**Returns** dataframe with imputations and their likelihoods in additional columns

**predict\_above\_precision** (*data\_frame*: *pandas.core.frame.DataFrame*, *precision\_threshold=0.95*) → dict

Returns the probabilities for each class, filtering out predictions below the precision threshold.

#### Parameters

- **data\_frame** – data frame
- **precision\_threshold** – don't predict if predicted class probability is below this precision threshold

**Returns** dict of {'column\_name': array}, array is a numpy array of shape samples-by-labels

**predict\_proba** (*data\_frame*: *pandas.core.frame.DataFrame*) → dict

Returns the probabilities for each class :param data\_frame: data frame :return: dict of {'column\_name': array}, array is a numpy array of shape samples-by-labels

**predict\_proba\_top\_k** (*data\_frame*: *pandas.core.frame.DataFrame*, *top\_k*: *int = 5*) → dict

Returns tuples of (label, probability) for the top\_k most likely predicted classes

#### Parameters

- **data\_frame** – pandas data frame
- **top\_k** – number of most likely predictions to return

**Returns** dict of {'column\_name': list} where list is a list of (label, probability) tuples

**save** ()

Saves model to disk, except mxnet module which is stored separately during fit

**ttransform** (*data\_frame*: *pandas.core.frame.DataFrame*) → dict

Imputes values given an mxnet iterator (see iterators) :param data\_frame: pandas data frame (pandas) :return: dict of {'column\_name': list} where list contains the string predictions

**transform\_and\_compute\_metrics** (*data\_frame: pandas.core.frame.DataFrame, metrics\_path=None*) → dict  
Returns predictions and metrics (average and per class)

**Parameters**

- **data\_frame** – data frame
- **metrics\_path** – if not None and exists, metrics are serialized as json to this path.

**Returns**

### 1.4.3 Column Encoders

Column Encoders: used for translating values of a table into numerical representation such that Featurizers can operate on them

**class** datawig.column\_encoders.**BowEncoder** (*input\_columns: Any, output\_column: str = None, max\_tokens: int = 262144, tokens: str = 'chars', ngram\_range: tuple = None, prefixed\_concatenation: bool = True*)

Bag-of-Words encoder for text data, using sklearn's HashingVectorizer

**Parameters**

- **input\_columns** – List[str] with column names to be used as input for this ColumnEncoder
- **output\_column** – Name of output field, used as field name in downstream MxNet iterator
- **max\_tokens** – Number of hash buckets (dimensionality of sparse ngram vector). default  $2^{18}$
- **tokens** – How to tokenize the input data, supports 'words' and 'chars'.
- **ngram\_range** – length of ngrams to use as features
- **prefixed\_concatenation** – whether or not to prefix values with column name before concat

**decode** (*col: pandas.core.series.Series*) → pandas.core.series.Series

Raises NotImplementedError, hashed bag-of-words cannot be decoded due to hash collisions

**Parameters** **token\_index\_sequence** –

**Returns**

**fit** (*data\_frame: pandas.core.frame.DataFrame*)

Does nothing, HashingVectorizers do not need to be fit.

**Parameters** **data\_frame** –

**Returns**

**is\_fitted** () → bool

Returns true if the column encoder does not require fitting (anymore or at all)

**Parameters** **self** –

**Returns** True if the encoder is fitted

**transform** (*data\_frame: pandas.core.frame.DataFrame*) → numpy.array

Transforms one or more string columns into Bag-of-words vectors, hashed into a max\_features dimensional feature space. Nans and missing values will be replaced by zero vectors.

**Parameters** **data\_frame** – pandas DataFrame with text columns

**Returns** numpy array (rows by max\_features)

**class** datawig.column\_encoders.**CategoricalEncoder** (*input\_columns: Any, output\_column: str = None, token\_to\_idx: Dict[str, int] = None, max\_tokens: int = 10000*)

Transforms categorical variable from string representation into number

#### Parameters

- **input\_columns** – List[str] with column names to be used as input for this ColumnEncoder
- **output\_column** – Name of output field, used as field name in downstream MxNet iterator
- **token\_to\_idx** – token to index mapping, 0 is reserved for missing tokens, 1 ... max\_tokens for most to least frequent tokens
- **max\_tokens** – maximum number of tokens

**decode** (*col: pandas.core.series.Series*) → pandas.core.series.Series

Decodes a pandas Series of token indices

**Parameters** **col** – pandas Series of token indices

**Returns** pandas Series of tokens

**decode\_token** (*token\_idx: int*) → str

Decodes a token index into a token

**Parameters** **token\_idx** – token index

**Returns** token

**fit** (*data\_frame: pandas.core.frame.DataFrame*)

Fits a CategoricalEncoder by extracting the value histogram of a column and capping it at max\_tokens. Issues warning if less than 100 values were observed.

**Parameters** **data\_frame** – pandas data frame

**is\_fitted** ()

Checks if ColumnEncoder (still) needs to be fitted to data

**Returns** True if the column encoder does not require fitting (anymore or at all)

**transform** (*data\_frame: pandas.core.frame.DataFrame*) → numpy.array

Transforms string column of pandas dataframe into categoricals

**Parameters** **data\_frame** – pandas data frame

**Returns** numpy array (rows by 1)

**static transform\_func\_categorical** (*col: pandas.core.series.Series, token\_to\_idx: Dict[str, int], missing\_token\_idx: int*) → Any

Transforms categorical values into their indices

#### Parameters

- **col** – pandas Series with categorical values

- **token\_to\_idx** – Dict[str, int] with mapping from token to token index
- **missing\_token\_idx** – index for missing symbol

### Returns

**class** datawig.column\_encoders.**ColumnEncoder** (*input\_columns: List[str], output\_column=None, output\_dim=1*)

Abstract super class of column encoders. Transforms value representation of columns (e.g. strings) into numerical representations to be fed into MxNet.

Options for ColumnEncoders are:

SequentialEncoder: for sequences of symbols (e.g. characters or words), BowEncoder: bag-of-word representation, as sparse vectors CategoricalEncoder: for categorical variables NumericalEncoder: for numerical values

### Parameters

- **input\_columns** – List[str] with column names to be used as input for this ColumnEncoder
- **output\_column** – Name of output field, used as field name in downstream MxNet iterator
- **output\_dim** – dimensionality of encoded column values (1 for categorical, vocabulary size for sequential and BoW)

**decode** (*col: pandas.core.series.Series*) → pandas.core.series.Series

Decodes a pandas Series of token indices

**Parameters** **col** – pandas Series of token indices

**Returns** pandas Series of tokens

**fit** (*data\_frame: pandas.core.frame.DataFrame*)

Fits a ColumnEncoder if needed (i.e. vocabulary/alphabet)

**Parameters** **data\_frame** – pandas DataFrame

**Returns**

**is\_fitted**()

Checks if ColumnEncoder (still) needs to be fitted to data

**Returns** True if the column encoder does not require fitting (anymore or at all)

**transform** (*data\_frame: pandas.core.frame.DataFrame*) → numpy.array

Transforms values in one or more columns of DataFrame into a numpy array that can be fed into a Featurizer

**Parameters** **data\_frame** –

**Returns** List of integers

**exception** datawig.column\_encoders.**NotFittedError**

Error thrown when unfitted encoder is used

**class** datawig.column\_encoders.**NumericalEncoder** (*input\_columns: Any, output\_column: str = None, normalize=True*)

Numerical encoder, concatenates columns in field\_names into one vector fills nans with the mean of a column

**Parameters**



- **input\_columns** – List[str] with column names to be used as input for this ColumnEncoder
- **output\_column** – Name of output field, used as field name in downstream MxNet iterator
- **normalize** – whether to normalize by the standard deviation or not, default True

**decode** (*col: pandas.core.series.Series*) → pandas.core.series.Series

Undoes the normalization, scales by scale and adds the mean

**Parameters** **col** – pandas Series (normalized)

**Returns** pandas Series (unnormalized)

**fit** (*data\_frame: pandas.core.frame.DataFrame*)

Does nothing or fits the normalizer, if normalization is specified

**Parameters** **data\_frame** – DataFrame with numerical columns specified when instantiating NumericalEncoder

**is\_fitted** ()

Returns true if the column encoder does not require fitting (anymore or at all)

**Parameters** **self** –

**Returns** True if the encoder is fitted

**transform** (*data\_frame: pandas.core.frame.DataFrame*) → numpy.array

Concatenates the numerical columns specified when instantiating the NumericalEncoder Normalizes features if specified in the NumericalEncoder

**Parameters** **data\_frame** – DataFrame with numerical columns specified in NumericalEncoder

**Returns** np.array with numerical features (rows by number of numerical columns)

**class** datawig.column\_encoders.**SequentialEncoder** (*input\_columns: Any, output\_column: str = None, token\_to\_idx: Dict[str, int] = None, max\_tokens: int = 1000, seq\_len: int = 500*)

Transforms sequence of characters into sequence of numbers

**Parameters**

- **input\_columns** – List[str] with column names to be used as input for this ColumnEncoder
- **output\_column** – Name of output field, used as field name in downstream MxNet iterator
- **token\_to\_idx** – token to index mapping 0 is reserved for missing tokens, 1 ... max\_tokens-1 for most to least frequent tokens
- **max\_tokens** – maximum number of tokens
- **seq\_len** – length of sequence, shorter sequences get padded to, longer sequences truncated at seq\_len symbols

**decode** (*col: pandas.core.series.Series*) → pandas.core.series.Series

Decodes a pandas Series of token indices

**Parameters** **col** – pandas Series of token index iterables

**Returns** pd.Series of strings

**decode\_seq** (*token\_index\_sequence: Iterable[int]*) → str

Decodes a sequence of token indices into a string

**Parameters** **token\_index\_sequence** – an iterable of token indices

**Returns** str the decoded string

**fit** (*data\_frame: pandas.core.frame.DataFrame*)

Fits a SequentialEncoder by extracting the character value histogram of a column and capping it at max\_tokens

**Parameters** **data\_frame** – pandas data frame

**is\_fitted** () → bool

Checks if ColumnEncoder (still) needs to be fitted to data

**Returns** True if the column encoder does not require fitting (anymore or at all)

**transform** (*data\_frame: pandas.core.frame.DataFrame*) → numpy.array

Transforms column of pandas dataframe into sequence of tokens

**Parameters** **data\_frame** – pandas DataFrame

**Returns** numpy array (rows by seq\_len)

**static transform\_func\_seq\_single** (*string: str, token\_to\_idx: Dict[str, int], seq\_len: int, missing\_token\_idx: int*) → List[int]

Transforms a single string into a sequence of token ids

**Parameters**

- **string** – a sequence of symbols as string
- **token\_to\_idx** – Dict[str, int] with mapping from token to token index
- **seq\_len** – length of sequence
- **missing\_token\_idx** – index for missing symbol

**Returns** List[int] with transformed values

**class** datawig.column\_encoders.**TfIdfEncoder** (*input\_columns: Any, output\_column: str = None, max\_tokens: int = 262144, tokens: str = 'chars', ngram\_range: tuple = None, prefixed\_concatenation: bool = True*)

TfIdf bag of word encoder for text data, using sklearn's TfidfVectorizer

**Parameters**

- **input\_columns** – List[str] with column names to be used as input for this ColumnEncoder
- **output\_column** – Name of output field, used as field name in downstream MxNet iterator
- **max\_tokens** – Number of feature buckets (dimensionality of sparse ngram vector). default 2\*\*18
- **tokens** – How to tokenize the input data, supports 'words' and 'chars'.
- **ngram\_range** – length of ngrams to use as features
- **prefixed\_concatenation** – whether or not to prefix values with column name before concat

**decode** (*col: pandas.core.series.Series*) → pandas.core.series.Series

Given a series of indices, decode it to input tokens

**Parameters** `col` –

**Returns** `pd.Series` of tokens

**fit** (*data\_frame: pandas.core.frame.DataFrame*)

**Parameters** `data_frame` –

**Returns**

**is\_fitted** () → `bool`

**Parameters** `self` –

**Returns** True if the encoder is fitted

**transform** (*data\_frame: pandas.core.frame.DataFrame*) → `numpy.array`

Transforms one or more string columns into Bag-of-words vectors.

**Parameters** `data_frame` – pandas DataFrame with text columns

**Returns** `numpy array` (rows by `max_features`)



**d**

`datawig.column_encoders`, 18  
`datawig.imputer`, 15  
`datawig.simple_imputer`, 10



**B**

BowEncoder (class in `datawig.column_encoders`), 18

**C**

calibrate() (`datawig.imputer.Imputer` method), 15

CategoricalEncoder (class in `datawig.column_encoders`), 19

check\_data\_types() (`datawig.simple_imputer.SimpleImputer` method), 11

check\_for\_label\_shift() (`datawig.simple_imputer.SimpleImputer` method), 11

ColumnEncoder (class in `datawig.column_encoders`), 20

complete() (`datawig.simple_imputer.SimpleImputer` static method), 11

**D**

`datawig.column_encoders` (module), 18

`datawig.imputer` (module), 15

`datawig.simple_imputer` (module), 10

decode() (`datawig.column_encoders.BowEncoder` method), 18

decode() (`datawig.column_encoders.CategoricalEncoder` method), 19

decode() (`datawig.column_encoders.ColumnEncoder` method), 20

decode() (`datawig.column_encoders.NumericalEncoder` method), 21

decode() (`datawig.column_encoders.SequentialEncoder` method), 21

decode() (`datawig.column_encoders.TfIdfEncoder` method), 22

decode\_seq() (`datawig.column_encoders.SequentialEncoder` method), 21

decode\_token() (`datawig.column_encoders.CategoricalEncoder` method), 20

**E**

explain() (`datawig.imputer.Imputer` method), 15

explain() (`datawig.simple_imputer.SimpleImputer` method), 12

explain\_instance() (`datawig.imputer.Imputer` method), 16

explain\_instance() (`datawig.simple_imputer.SimpleImputer` method), 12

**F**

fit() (`datawig.column_encoders.BowEncoder` method), 18

fit() (`datawig.column_encoders.CategoricalEncoder` method), 19

fit() (`datawig.column_encoders.ColumnEncoder` method), 20

fit() (`datawig.column_encoders.NumericalEncoder` method), 21

fit() (`datawig.column_encoders.SequentialEncoder` method), 22

fit() (`datawig.column_encoders.TfIdfEncoder` method), 23

fit() (`datawig.imputer.Imputer` method), 16

fit() (`datawig.simple_imputer.SimpleImputer` method), 12

fit\_hpo() (`datawig.simple_imputer.SimpleImputer` method), 13

Imputer (class in `datawig.imputer`), 15

is\_fitted() (`datawig.column_encoders.BowEncoder` method), 18

is\_fitted() (`datawig.column_encoders.CategoricalEncoder` method), 19

is\_fitted() (`datawig.column_encoders.ColumnEncoder` method), 20

is\_fitted() (`datawig.column_encoders.NumericalEncoder` method), 21

is\_fitted() (*datawig.column\_encoders.SequentialEncoder method*), 22  
 is\_fitted() (*datawig.column\_encoders.TfIdfEncoder method*), 23  
 is\_fitted() (*datawig.imputer.Imputer method*), 17  
 transform\_and\_compute\_metrics() (*datawig.imputer.Imputer method*), 17  
 transform\_func\_categorical() (*datawig.column\_encoders.CategoricalEncoder static method*), 19

## L

load() (*datawig.imputer.Imputer static method*), 16  
 load() (*datawig.simple\_imputer.SimpleImputer static method*), 14  
 load\_hpo\_model() (*datawig.simple\_imputer.SimpleImputer method*), 14  
 load\_metrics() (*datawig.simple\_imputer.SimpleImputer method*), 14  
 transform\_func\_seq\_single() (*datawig.column\_encoders.SequentialEncoder static method*), 22

## N

NotFittedError, 20  
 NumericalEncoder (class in *datawig.column\_encoders*), 20

## P

predict() (*datawig.imputer.Imputer method*), 16  
 predict() (*datawig.simple\_imputer.SimpleImputer method*), 14  
 predict\_above\_precision() (*datawig.imputer.Imputer method*), 17  
 predict\_proba() (*datawig.imputer.Imputer method*), 17  
 predict\_proba\_top\_k() (*datawig.imputer.Imputer method*), 17

## S

save() (*datawig.imputer.Imputer method*), 17  
 save() (*datawig.simple\_imputer.SimpleImputer method*), 15  
 SequentialEncoder (class in *datawig.column\_encoders*), 21  
 SimpleImputer (class in *datawig.simple\_imputer*), 10

## T

TfIdfEncoder (class in *datawig.column\_encoders*), 22  
 transform() (*datawig.column\_encoders.BowEncoder method*), 18  
 transform() (*datawig.column\_encoders.CategoricalEncoder method*), 19  
 transform() (*datawig.column\_encoders.ColumnEncoder method*), 20  
 transform() (*datawig.column\_encoders.NumericalEncoder method*), 21  
 transform() (*datawig.column\_encoders.SequentialEncoder method*), 22  
 transform() (*datawig.column\_encoders.TfIdfEncoder method*), 23