
Gretel Synthetics

Gretel.ai

Apr 02, 2024

CONTENTS:

1	Documentation	3
2	Try it out now!	5
3	Getting Started	7
3.1	Dependency Requirements	7
4	Timeseries DGAN Overview	9
5	ACTGAN Overview	11
6	LSTM Overview	13
6.1	Simple Mode	13
6.2	DataFrame Mode	13
6.3	Components	13
6.4	Utilities	14
6.5	Differential Privacy	14
7	Modules	15
7.1	Config	15
7.2	Tokenizers	19
7.3	Train	22
7.4	Generate	22
7.5	Batch	25
7.6	Utils	32
7.7	Timeseries DGAN	35
7.8	ACTGAN	42
8	Indices and tables	49
	Python Module Index	51
	Index	53

gretel

DOCUMENTATION

- [Get started with gretel-synthetics](#)
- [Configuration](#)
- [Train your model](#)
- [Generate synthetic records](#)

TRY IT OUT NOW!

If you want to quickly discover gretel-synthetics, simply click the button below and follow the tutorials!

Check out additional examples [here](#).

GETTING STARTED

This section will guide you through installation of `gretel-synthetics` and dependencies that are not directly installed by the Python package manager.

3.1 Dependency Requirements

By default, we do not install certain core requirements, the following dependencies should be installed *external to the installation* of `gretel-synthetics`, depending on which model(s) you plan to use.

- Tensorflow: Used by the LSTM model, we recommend version 2.11.x
- Torch: Used by Timeseries DGAN and ACTGAN (for ACTGAN, Torch is installed by SDV), we recommend version 2.0
- SDV (Synthetic Data Vault): Used by ACTGAN, we recommend version 0.17.x

These dependencies can be installed by doing the following:

```
pip install tensorflow==2.11 # for LSTM
pip install sdv<0.18 # for ACTGAN
pip install torch==2.0 # for Timeseries DGAN
```

To install the actual `gretel-synthetics` package, first clone the repo and then...

```
pip install -U .
```

or

```
pip install gretel-synthetics
```

then...

```
$ pip install jupyter
$ jupyter notebook
```

When the UI launches in your browser, navigate to `examples/synthetic_records.ipynb` and get generating!

If you want to install `gretel-synthetics` locally and use a GPU (recommended):

1. Create a virtual environment (e.g. using conda)

```
$ conda create --name tf python=3.9
```

1. Activate the virtual environment

```
$ conda activate tf
```

1. Run the setup script `./setup-utils/setup-gretel-synthetics-tensorflow24-with-gpu.sh`

The last step will install all the necessary software packages for GPU usage, tensorflow=2.8 and gretel-synthetics. Note that this script works only for Ubuntu 18.04. You might need to modify it for other OS versions.

TIMESERIES DGAN OVERVIEW

The [timeseries DGAN module](#) contains a PyTorch implementation of a DoppelGANger model that is optimized for timeseries data. Similar to tensorflow, you will need to manually install pytorch:

```
pip install torch==1.13.1
```

This [notebook](#) shows basic usage on a small data set of smart home sensor readings.

ACTGAN OVERVIEW

ACTGAN (Anyway CTGAN) is an extension of the popular [CTGAN implementation](#) that provides some additional functionality to improve memory usage, autodetection and transformation of columns, and more.

To use this model, you will need to manually install SDV:

```
pip install sdv<0.18
```

Keep in mind that this will also install several dependencies like PyTorch that SDV relies on, which may conflict with PyTorch versions installed for use with other models like Timeseries DGAN.

The ACTGAN interface is a superset of the CTGAN interface. To see the additional features, please take a look at the ACTGAN demo notebook in the `examples` directory of this repo.

LSTM OVERVIEW

This package allows developers to quickly get immersed with synthetic data generation through the use of neural networks. The more complex pieces of working with libraries like Tensorflow and differential privacy are bundled into friendly Python classes and functions. There are two high level modes that can be utilized.

6.1 Simple Mode

The simple mode will train line-per-line on an input file of text. When generating data, the generator will yield a custom object that can be used a variety of different ways based on your use case. [This notebook](#) demonstrates this mode.

6.2 DataFrame Mode

This library supports CSV / DataFrames natively using the DataFrame “batch” mode. This module provided a wrapper around our simple mode that is geared for working with tabular data. Additionally, it is capable of handling a high number of columns by breaking the input DataFrame up into “batches” of columns and training a model on each batch. [This notebook](#) shows an overview of using this library with DataFrames natively.

6.3 Components

There are four primary components to be aware of when using this library.

1. **Configurations.** Configurations are classes that are specific to an underlying ML engine used to train and generate data. An example would be using `TensorFlowConfig` to create all the necessary parameters to train a model based on TF. `LocalConfig` is aliased to `TensorFlowConfig` for backwards compatibility with older versions of the library. A model is saved to a designated directory, which can optionally be archived and utilized later.
2. **Tokenizers.** Tokenizers convert input text into integer based IDs that are used by the underlying ML engine. These tokenizers can be created and sent to the training input. This is optional, and if no specific tokenizer is specified then a default one will be used. You can find [an example](#) here that uses a simple char-by-char tokenizer to build a model from an input CSV. When training in a non-differentially private mode, we suggest using the default `SentencePiece` tokenizer, an unsupervised tokenizer that learns subword units (e.g., **byte-pair-encoding (BPE)** [Sennrich et al.]) and **unigram language model** [Kudo.]) for faster training and increased accuracy of the synthetic model.
3. **Training.** Training a model combines the configuration and tokenizer and builds a model, which is stored in the designated directory, that can be used to generate new records.

4. Generation. Once a model is trained, any number of new lines or records can be generated. Optionally, a record validator can be provided to ensure that the generated data meets any constraints that are necessary. See our notebooks for examples on validators.

6.4 Utilities

In addition to the four primary components, the `gretel-synthetics` package also ships with a set of utilities that are helpful for training advanced synthetics models and evaluating synthetic datasets.

Some of this functionality carries large dependencies, so they are shipped as an extra called `utils`. To install these dependencies, you may run

```
pip install gretel-synthetics[utils]
```

For additional details, please refer to the [Utility module API docs](#).

6.5 Differential Privacy

Differential privacy support for our TensorFlow mode is built on the great work being done by the Google TF team and their [TensorFlow Privacy library](#).

When utilizing DP, we currently recommend using the character tokenizer as it will only create a vocabulary of single tokens and removes the risk of sensitive data being memorized as actual tokens that can be replayed during generation.

There are also a few configuration options that are notable such as:

- `predict_batch_size` should be set to 1
- `dp` should be enabled
- `learning_rate`, `dp_noise_multiplier`, `dp_l2_norm_clip`, and `dp_microbatches` can be adjusted to achieve various epsilon values.
- `reset_states` should be disabled

Please see our [example Notebook](#) for training a DP model based on the [Netflix Prize](#) dataset.

MODULES

7.1 Config

This module provides a set of dataclasses that can be used to hold all necessary configuration parameters for training a model and generating data.

For example usage please see our Jupyter Notebooks.

```
class gretel_synthetics.config.BaseConfig(input_data_path: str | None = None, validation_split: bool =
    True, checkpoint_dir: str | None = None, training_data_path:
    str | None = None, field_delimiter: str | None = None,
    field_delimiter_token: str = '<d>', model_type: str | None =
    None, max_lines: int = 0, overwrite: bool = False,
    epoch_callback: Callable | None = None,
    max_training_time_seconds: int | None = None, vocab_size:
    int = 20000, character_coverage: float = 1.0,
    pretrain_sentence_count: int = 1000000, max_line_len: int =
    2048)
```

This class should not be used directly, engine specific classes should derived from this class.

as_dict()

Serialize the config attrs to a dict

checkpoint_dir: str = None

Directory where model data will be stored, user provided.

epoch_callback: Callable | None = None

Callback to be invoked at the end of each epoch. It will be invoked with an EpochState instance as its only parameter. NOTE that the callback is deleted when save_model_params is called, we do not attempt to serialize it to JSON.

field_delimiter: str | None = None

If the input data is structured, you may specify a field delimiter which can be used to split the generated text into a list of strings. For more detail please see the GenText class in the generate.py module.

field_delimiter_token: str = '<d>'

Depending on the tokenizer used, a special token can be used to represent characters. For tokenizers, like SentencePiece that support this, we will replace the field delimiter char with this token to provide better learning and generation. If the tokenizer used does not support custom tokens, this value will be ignored

abstract get_generator_class() → None

This must be implemented by all specific configs. It should return the class that should be used as the Generator for creating records.

abstract get_training_callable() → Callable

This must be implemented by all specific configs. It should return a callable that should be used as the entrypoint for training a model.

gpu_check()

Optionally do a GPU check and warn if a GPU is not available, if not overridden, do nothing

input_data_path: str = None

Path to raw training data, user provided.

max_lines: int = 0

The maximum number of lines to utilize from the raw input data.

max_training_time_seconds: int | None = None

If set, training will cease after the number of seconds specified elapses. This timeout will be evaluated after each epoch.

model_type: str = None

A string version of the model config class. This is used to keep track of what underlying engine was used when writing the config to a file. This will be automatically updated during construction.

overwrite: bool = False

Set to True to automatically overwrite previously saved model checkpoints. If False, the trainer will generate an error if checkpoints exist in the model directory. Default is False.

training_data_path: str = None

Where annotated and tokenized training data will be stored. This attr will be modified during construction.

validation_split: bool = True

Use a fraction of the training data as validation data. Use of a validation set is recommended as it helps prevent over-fitting and memorization. When enabled, 20% of data will be used for validation.

**gretel_synthetics.config.CONFIG_MAP = {'TensorFlowConfig': <class
'gretel_synthetics.config.TensorFlowConfig'>}**

A mapping of configuration subclass string names to their actual classes. This can be used to re-instantiate a config from a serialized state.

gretel_synthetics.config.LocalConfig

alias of *TensorFlowConfig*

```

class gretel_synthetics.config.TensorFlowConfig(input_data_path: str | None = None, validation_split:
    bool = True, checkpoint_dir: str | None = None,
    training_data_path: str | None = None,
    field_delimiter: str | None = None,
    field_delimiter_token: str = '<d>', model_type: str |
    None = None, max_lines: int = 0, overwrite: bool =
    False, epoch_callback: Callable | None = None,
    max_training_time_seconds: int | None = None,
    vocab_size: int = 20000, character_coverage: float =
    1.0, pretrain_sentence_count: int = 1000000,
    max_line_len: int = 2048, epochs: int = 100,
    early_stopping: bool = True,
    early_stopping_patience: int = 5, best_model_metric:
    str | None = None, early_stopping_min_delta: float =
    0.001, batch_size: int = 64, buffer_size: int = 10000,
    seq_length: int = 100, embedding_dim: int = 256,
    rnn_units: int = 256, learning_rate: float = 0.01,
    dropout_rate: float = 0.2, rnn_initializer: str =
    'glorot_uniform', dp: bool = False,
    dp_noise_multiplier: float = 0.1, dp_l2_norm_clip:
    float = 3.0, dp_microbatches: int = 1, gen_temp: float
    = 1.0, gen_chars: int = 0, gen_lines: int = 1000,
    predict_batch_size: int = 64, reset_states: bool =
    True, save_all_checkpoints: bool = False,
    save_best_model: bool = True)

```

TensorFlow config that contains all of the main parameters for training a model and generating data.

Parameters

- **epochs** (*optional*) – Number of epochs to train the model. An epoch is an iteration over the entire training set provided. For production use cases, 15-50 epochs are recommended. The default is 100 and is intentionally set extra high. By default, **early_stopping** is also enabled and will stop training epochs once the model is no longer improving.
- **early_stopping** (*optional*) – deduce when the model is no longer improving and terminating training.
- **early_stopping_patience** (*optional*) – in the model. After this number of epochs, training will terminate.
- **best_model_metric** (*optional*) – The metric to use to track when a model is no longer improving. Alternative options are “val_acc” or “acc”. A error will be raised if a valid value is not specified.
- **early_stopping_min_delta** (*optional*) – as an improvement, i.e. an absolute change of less than min_delta will count as no improvement.
- **batch_size** (*optional*) – Number of samples per gradient update. Using larger batch sizes can help make more efficient use of CPU/GPU parallelization, at the cost of memory. If unspecified, batch_size will default to 64.
- **buffer_size** (*optional*) – Buffer size which is used to shuffle elements during training. Default size is 10000.
- **seq_length** (*optional*) – The maximum length sentence we want for a single training input in characters. Note that this setting is different than max_line_length, as seq_length simply affects the length of the training examples passed to the neural network to predict the next token. Default size is 100.

- **embedding_dim** (*optional*) – Vector size for the lookup table used in the neural network Embedding layer that maps the numbers of each character. Default size is 256.
- **rnn_units** (*optional*) – Positive integer, dimensionality of the output space for LSTM layers. Default size is 256.
- **dropout_rate** (*optional*) – Float between 0 and 1. Fraction of the units to drop for the linear transformation of the inputs. Using a dropout can help to prevent overfitting by ignoring randomly selected neurons during training. 0.2 (20%) is often used as a good compromise between retaining model accuracy and preventing overfitting. Default is 0.2.
- **rnn_initializer** (*optional*) – Initializer for the kernel weights matrix, used for the linear transformation of the inputs. Default is `glorot_transform`.
- **dp** (*optional*) – If True, train model with differential privacy enabled. This setting provides assurances that the models will encode general patterns in data rather than facts about specific training examples. These additional guarantees can usefully strengthen the protections offered for sensitive data and content, at a small loss in model accuracy and synthetic data quality. The differential privacy epsilon and delta values will be printed when training completes. Default is False.
- **learning_rate** (*optional*) – The higher the learning rate, the more that each update during training matters. Note: When training with differential privacy enabled, if the updates are noisy (such as when the additive noise is large compared to the clipping threshold), a low learning rate may help with training. Default is `0.01`.
- **dp_noise_multiplier** (*optional*) – The amount of noise sampled and added to gradients during training. Generally, more noise results in better privacy, at the expense of model accuracy. Default is `0.1`.
- **dp_l2_norm_clip** (*optional*) – The maximum Euclidean (L2) norm of each gradient is applied to update model parameters. This hyperparameter bounds the optimizer's sensitivity to individual training points. Default is `3.0`.
- **dp_microbatches** (*optional*) – Each batch of data is split into smaller units called micro-batches. Computational overhead can be reduced by increasing the size of micro-batches to include more than one training example. The number of micro-batches should divide evenly into the overall `batch_size`. Default is 1.
- **gen_temp** (*optional*) – Controls the randomness of predictions by scaling the logits before applying softmax. Low temperatures result in more predictable text. Higher temperatures result in more surprising text. Experiment to find the best setting. Default is `1.0`.
- **gen_chars** (*optional*) – Maximum number of characters to generate per line. Default is 0 (no limit).
- **gen_lines** (*optional*) – Maximum number of text lines to generate. This function is used by `generate_text` and the optional `line_validator` to make sure that all lines created by the model pass validation. Default is `1000`.
- **predict_batch_size** (*optional*) – How many words to generate in parallel. Higher values may result in increased throughput. The default of 64 should provide reasonable performance for most users.
- **reset_states** (*optional*) – Reset RNN model states between each record created guarantees more consistent record creation over time, at the expense of model accuracy. Default is True.
- **save_all_checkpoints** (*optional*) – which can be useful for optimal model selection. Set to False to save only the latest checkpoint. Default is True.

- **save_best_model** (optional). Defaults to True. Track the best version of the model (checkpoint) – If **save_all_checkpoints** is disabled, then the saved model will be overwritten by newer ones only if they are better.

get_generator_class()

This must be implemented by all specific configs. It should return the class that should be used as the Generator for creating records.

get_training_callable()

This must be implemented by all specific configs. It should return a callable that should be used as the entrypoint for training a model.

gpu_check()

Optionally do a GPU check and warn if a GPU is not available, if not overridden, do nothing

`gretel_synthetics.config.config_from_model_dir(model_dir: str) → BaseConfig`

Factory that will take a known directory of a model and return a class instance for that config. We automatically try and detect the correct BaseConfig sub-class to use based on the saved model params.

If there is no `model_type` param in the saved config, we assume that the model was saved using an earlier version of the package and will instantiate a TensorFlowConfig

7.2 Tokenizers

Interface definitions for tokenizers. The classes in the module are segmented into two abstract types: Trainers and Tokenizers. They are kept separate because the parameters used to train a tokenizer are not necessarily loaded back in and utilized by a trained tokenizer. While its more explicit to utilize two types of classes, it also removes any ambiguity in which methods are able to be used based on training or tokenizing.

Trainers require a specific configuration to be provided. Based on the configuration received, the tokenizer trainers will create the actual training data file that will be used by the downstream training process. In this respect, utilizing at least one of these tokenizers is required for training since it is the tokenizers responsibility to create the final training data to be used.

The general process that is followed when using these tokenizers is:

Create a trainer instance, with desired parameters, including providing the config as a required param.

Call the `annotate_data` for your tokenizer trainer. What is important to note here is that this method actually iterates the input data line by line, and does any special processing, then writes a new data file that will be used for actual training. This new data file is written to the model directory.

Call the `train` method, which will create your tokenization model and save it to the model directory.

Now you will use the `load()` class method from an actual tokenizer class to load that trained model in and now you can use it on input data.

class `gretel_synthetics.tokenizers.Base`

High level base class for shared class attrs and validation. Should not be used directly.

class `gretel_synthetics.tokenizers.BaseTokenizer(model_data: Any, model_dir: str)`

Base class for loading a tokenizer from disk. Should not be used directly.

decode_from_ids(ids: List[int]) → str

Given a list of token IDs, convert it to a single string that would be the original string it was.

Note: We automatically call a method that can optionally restore any special reserved tokens back to their original values (such as field delimiter values, etc)

encode_to_ids(*data: str*) → List[int]

Given an input string, convert it to a list of token IDs

abstract classmethod load(*model_dir: str*)

Given a directory to a model, load the specific tokenizer model into an instance. Subclasses should implement this logic specific to how they need to load a model back in

abstract property total_vocab_size

Return the total count of unique tokens in the vocab, specific to the underlying tokenizer to be used.

class gretel_synthetics.tokenizers.**BaseTokenizerTrainer**(* , *config: None, vocab_size: int | None = None*)

Base class for training tokenizers. Should not be used directly.

annotate_data() → Iterator[str]

This should be called `_before_training` as it is required to have the annotated training data created in the model directory.

Read in the configurations raw input data path, and create a file I/O pipeline where each line of the input data path can optionally route through an annotation function and then we will write each raw line out into a training data file as specified by the config.

config: None

A subclass instance of `BaseConfig`. This will be used to find the input data for tokenization

data_iterator() → Iterator[str]

Create a generator that will iterate each line of the training data that was created during the annotation step. Synthetic model trainers will most likely need to iterate this to process each line of the annotated training data.

num_lines: int = 0

The number of lines that were processed after `create_annotated_training_data` is called

train()

Train a tokenizer and save the tokenizer settings to a file located in the model directory specified by the config object

vocab_size: int

The max size of the vocab (tokens) to be extracted from the input dataset.

class gretel_synthetics.tokenizers.**CharTokenizer**(*model_data: Any, model_dir: str*)

Load a simple character tokenizer from disk to conduct encoding and decoding operations

classmethod load(*model_dir: str*)

Create an instance of this tokenizer.

Parameters

model_dir – The path to the model directory

property total_vocab_size

Get the number of unique characters (tokens)


```
class gretel_synthetics.tokenizers.CharTokenizerTrainer(*, config: None, vocab_size: int | None = None)
```

Train a simple tokenizer that maps every single character to a unique ID. If `vocab_size` is not specified, the learned vocab size will be the number of unique characters in the training dataset.

Parameters

vocab_size – Max number of tokens (chars) to map to tokens.

```
class gretel_synthetics.tokenizers.SentencePieceColumnTokenizer(sp: SentencePieceProcessor, model_dir: str)
```

```
class gretel_synthetics.tokenizers.SentencePieceColumnTokenizerTrainer(col_pattern: str = '<col{>}', **kwargs)
```

```
class gretel_synthetics.tokenizers.SentencePieceTokenizer(model_data: Any, model_dir: str)
```

Load a SentencePiece tokenizer from disk so encoding / decoding can be done

```
classmethod load(model_dir: str)
```

Load a SentencePiece tokenizer from a model directory.

Parameters

model_dir – The model directory.

property total_vocab_size

The number of unique tokens in the model

```
class gretel_synthetics.tokenizers.SentencePieceTokenizerTrainer(*, character_coverage: float = 1.0, pretrain_sentence_count: int = 1000000, max_line_len: int = 2048, **kwargs)
```

Train a tokenizer using Google SentencePiece.

character_coverage: float

The amount of characters covered by the model. Unknown characters will be replaced with the <unk> tag. Good defaults are 0.995 for languages with rich character sets like Japanese or Chinese, and 1.0 for other languages or machine data. Default is 1.0.

max_line_len: int

Maximum line length for input training data. Any lines longer than this length will be ignored. Default is 2048.

pretrain_sentence_count: int

The number of lines `spm_train` first loads. Remaining lines are simply discarded. Since `spm_train` loads entire corpus into memory, this size will depend on the memory size of the machine. It also affects training time. Default is 1000000.

vocab_size: int

Pre-determined maximum vocabulary size prior to neural model training, based on subword units including byte-pair-encoding (BPE) and unigram language model, with the extension of direct training from raw sentences. We generally recommend using a large vocabulary size of 20,000 to 50,000. Default is 20000.

```
exception gretel_synthetics.tokenizers.TokenizerError
```

```
exception gretel_synthetics.tokenizers.VocabSizeTooSmall
```

Error that is raised when the `vocab_size` is too small for the given data. This happens, when the `vocab_size` is set to a value that is smaller than the number of required characters.

`gretel_synthetics.tokenizers.tokenizer_from_model_dir(model_dir: str) → BaseTokenizer`

A factory function that will return a tokenizer instance that can be used for encoding / decoding data. It will try to automatically infer what type of class to use based on the stored tokenizer params in the provided model directory.

If no specific tokenizer type is found, we assume that we are restoring a SentencePiece tokenizer because the model is from a version $\leq 0.14.x$

Parameters

model_dir – A directory that holds synthetic model data.

7.3 Train

Train models for creating synthetic data. This module is the primary entrypoint for creating a model. It depends on having created an engine specific configuration and optionally a tokenizer to be used.

class `gretel_synthetics.train.EpochState`(*epoch: int, accuracy: float | None = None, loss: float | None = None, val_accuracy: float | None = None, val_loss: float | None = None, batch: int | None = None, epsilon: float | None = None, delta: float | None = None*)

Training state passed to the epoch callback on BaseConfig at the end of each epoch.

class `gretel_synthetics.train.TrainingParams`(*tokenizer_trainer: None, tokenizer: None, config: None*)

A structure that is created and passed into the engine-specific training entrypoint. All engine-specific training entrypoints should expect to receive this object and process accordingly.

`gretel_synthetics.train.train`(*store: None, tokenizer_trainer: None = None*)

Train a Synthetic Model. This is a facade entrypoint that implements the engine specific training operation based on the provided configuration.

Parameters

- **store** – A subclass instance of BaseConfig. This config is responsible for providing the actual training entrypoint for a specific training routine.
- **tokenizer_trainer** – An optional subclass instance of a BaseTokenizerTrainer. If provided this tokenizer will be used to pre-process and create an annotated dataset for training. If not provided a default tokenizer will be used.

`gretel_synthetics.train.train_rnn`(*store: None*)

Facade to support backwards compatibility for $\leq 0.14.x$ versions.

7.4 Generate

Abstract module for generating data. The `generate_text` function is the primary entrypoint for creating text.

class `gretel_synthetics.generate.BaseGenerator`

Do not use directly.

Specific generation modules should have a subclass of this ABC that implements the core logic for generating data

class `gretel_synthetics.generate.GenText`(*valid: bool = None, text: str = None, explain: str = None, delimiter: str = None*)

`gretel_synthetics.generate.PredString`

alias of `pred_string`

class `gretel_synthetics.generate.SeedingGenerator`(*config: None, *, seed_list: List[str], line_validator: Callable | None = None, max_invalid: int = 1000*)

A single threaded line / text generator that is specifically for using with a list of seeds. This also exposes the `Settings` class back to the caller so the actual seed list can be directly accessed, which controls the underlying progression of the main text generator.

This is useful when you need to manipulate the actual seed list as data is being generated.

class `gretel_synthetics.generate.Settings`(*config: None, start_string: str | List[str] | None = None, multi_seed: bool = False, line_validator: Callable | None = None, max_invalid: int = 1000, tokenizer: BaseTokenizer | None = None, generator: BaseGenerator | None = None*)

Do not use directly.

Arguments for a generator generating lines of text.

This class contains basic settings for a generation process. It is separated from the `Generator` class for ensuring reliable serializability without an excess amount of code tied to it.

This class also will take a provided start string and validate that it can be utilized for text generation. If the `start_string` is something other than the default, we have to do a couple things:

- 1) If the config utilizes a field delimiter, the `start_string` MUST end with that delimiter
- 2) Convert the user-facing delim char into the special delim token specified in the config

class `gretel_synthetics.generate.gen_text`(*valid: bool | None = None, text: str | None = None, explain: str | None = None, delimiter: str | None = None*)

A record that is yielded from the `Generator.generate_next` generator.

valid

True, False, or None. If the line passed a validation function, then this will be `True`. If the validation function raised an exception then this will be automatically set to `False`. If no validation function is used, then this value will be `None`.

Type
bool

text

The actual record as a string

Type
str

explain

A string that describes why a record failed validation. This is the string representation of the `Exception` that is thrown in a validation function. This will only be set if validation fails, otherwise will be `None`.

Type
str

delimiter

If the generated text are column/field based records. This will hold the delimiter used to separate the fields from each other.

Type
str

as_dict() → dict

Serialize the generated record to a dictionary

values_as_list() → List[str] | None

Attempt to split the generated text on the provided delimiter

Returns

A list of values that are separated by the object's delimiter or None if there is no delimiter in the text

`gretel_synthetics.generate.generate_text`(*config: None, start_string: str | List[str] | None = None, line_validator: Callable | None = None, max_invalid: int = 1000, num_lines: int | None = None, parallelism: int = 0*) → Iterator[[GenText](#)]

A generator that will load a model and start creating records.

Parameters

- **config** – A configuration object, which you must have created previously
- **start_string** – A prefix string that is used to seed the record generation. By default we use a newline, but you may substitute any initial value here which will influence how the generator predicts what to generate. If you are working with a field delimiter, and you want to seed more than one column value, then you **MUST** utilize the field delimiter specified in your config. An example would be “foo,bar,baz,”. Also, if using a field delimiter, the string **MUST** end with the delimiter value.

Note: This param may also be a list of prefixes. If this is provided, then the generator will attempt to create exactly 1 record for each seed in the list. The `num_lines` param will be implicitly set to the size of the list and this number of records will be created at a 1:1 ratio between prefix strings and valid records.

- **line_validator** – An optional callback validator function that will take the raw string value from the generator as a single argument. This validator can execute arbitrary code with the raw string value. The validator function may return a bool to indicate line validity. This boolean value will be set on the yielded `gen_text` object. Additionally, if the validator throws an exception, the `gen_text` object will be set with a failed validation. If the validator returns None, we will assume successful validation.
- **max_invalid** – If using a `line_validator`, this is the maximum number of invalid lines to generate. If the number of invalid lines exceeds this value a `RuntimeError` will be raised.
- **num_lines** – If not None, this will override the `gen_lines` value that is provided in the config. .. note:

If `start_string` is a list, this value will be set to the length of that list and any other values for the param are ignored.

- **parallelism** – The number of concurrent workers to use. 1 (the default) disables parallelization, while a non-positive value means “number of CPUs + x” (i.e., use 0 for using as many workers as there are CPUs). A floating-point value is interpreted as a fraction of the available CPUs, rounded down.

Simple validator example:

```
def my_validator(raw_line: str):
    parts = raw_line.split(',')
    if len(parts) != 5:
        raise Exception('record does not have 5 fields')
```

Note: `gen_lines` from the `config` is important for this function. If a line validator is not provided, each line will count towards the number of total generated lines. When the total lines generated is \geq `gen_lines` we stop. If a line validator is provided, only *valid* lines will count towards the total number of lines generated. When the total number of valid lines generated is \geq `gen_lines`, we stop.

Note: `gen_chars`, controls the possible maximum number of characters a single generated line can have. If a newline character has not been generated before reaching this number, then the line will be returned. For example if `gen_chars` is 180 and a newline has not been generated, once 180 chars have been created, the line will be returned no matter what. As a note, if this value is 0, then each line will generate until a newline is observed.

Yields

A `GenText` object for each record that is generated. The generator will stop after the max number of lines is reached (based on your config).

Raises

A `RuntimeError` if the `max_invalid` number of lines is generated –

7.5 Batch

This module allows automatic splitting of a `DataFrame` into smaller `DataFrames` (by clusters of columns) and doing model training and text generation on each sub-DF independently.

Then we can concat each sub-DF back into one final synthetic dataset.

For example usage, please see our Jupyter Notebook.

```
class gretel_synthetics.batch.Batch(checkpoint_dir: str, input_data_path: str, headers: List[str], config:
    TensorFlowConfig, gen_data_count: int = 0)
```

A representation of a synthetic data workflow. It should not be used directly. This object is created automatically by the primary batch handler, such as `DataFrameBatch`. This class holds all of the necessary information for training, data generation and `DataFrame` re-assembly.

add_valid_data(data: `GenText`)

Take a `gen_text` object and add the generated line to the generated data stream

get_validator()

If a custom validator is set, we return that. Otherwise, we return the built-in validator, which simply checks if a generated line has the right number of values based on the number of headers for this batch.

This at least makes sure the resulting `DataFrame` will be the right shape

load_validator_from_file()

Load a saved validation object if it exists

reset_gen_data()

Reset all objects that accumulate or track synthetic data generation

set_validator(*fn: Callable, save=True*)

Assign a validation callable to this batch. Optionally pickling and saving the validator for loading later

property synthetic_df: DataFrame

Get a DataFrame constructed from the generated lines

```
class gretel_synthetics.batch.DataFrameBatch(*, df: DataFrame | None = None, batch_size: int = 15,
                                             batch_headers: List[List[str]] | None = None, config: dict
                                             | BaseConfig | None = None, tokenizer:
                                             BaseTokenizerTrainer | None = None, mode: str = 'write',
                                             checkpoint_dir: str | None = None, validate_model: bool
                                             = True)
```

Create a multi-batch trainer / generator. When created, the directory structure to store models and training data will automatically be created. The directory structure will be created under the “checkpoint_dir” location provided in the config template. There will be one directory per batch, where each directory will be called “batch_N” where N is the batch number, starting from 0.

Training and generating can happen per-batch or we can loop over all batches to do both train / generation functions.

Example

When creating this object, you must explicitly create the training data from the input DataFrame before training models:

```
my_batch = DataFrameBatch(df=my_df, config=my_config)
my_batch.create_training_data()
my_batch.train_all_batches()
```

Parameters

- **df** – The input, source DataFrame
- **batch_size** – If **batch_headers** is not provided we automatically break up the number of columns in the source DataFrame into batches of N columns.
- **batch_headers** – A list of lists of strings can be provided which will control the number of batches. The number of inner lists is the number of batches, and each inner list represents the columns that belong to that batch
- **config** – A template training config to use, this will be used as kwargs for each Batch’s synthetic configuration. This may also be a subclass of **BaseConfig**. If this is used, you can set the **input_data_path** param to the constant **PATH_HOLDER** as it does not really matter
- **tokenizer_class** – An optional **BaseTokenizerTrainer** subclass. If not provided the default tokenizer will be used for the underlying ML engine.

Note: When providing a config, the source of training data is not necessary, only the **checkpoint_dir** is needed. Each batch will control its input training data path after it creates the training dataset.

batch_size: int

The max number of columns allowed for a single DF batch

batch_to_df(*batch_idx: int*) → DataFrame

Extract a synthetic data DataFrame from a single batch.

Parameters

batch_idx – The batch number

Returns

A DataFrame with synthetic data

batches: Dict[int, Batch]

A mapping of Batch objects to a batch number. The batch number (key) increments from 0..N where N is the number of batches being used.

batches_to_df() → DataFrame

Convert all batches to a single synthetic data DataFrame.

Returns

A single DataFrame that is the concatenation of all the batch DataFrames.

config: dict | BaseConfig

The template config that will be used for all batches. If a dict is provided we default to a TensorFlowConfig.

create_training_data()

Split the original DataFrame into N smaller DataFrames. Each smaller DataFrame will have the same number of rows, but a subset of the columns from the original DataFrame.

This method iterates over each Batch object and assigns a smaller training DataFrame to the `training_df` attribute of the object.

Finally, a training CSV is written to disk in the specific batch directory

generate_all_batch_lines(*max_invalid=1000, raise_on_failed_batch: bool = False, num_lines: int | None = None, seed_fields: dict | List[dict] | None = None, parallelism: int = 0*) → Dict[int, GenerationSummary]

Generate synthetic lines for all batches. Lines for each batch are added to the individual Batch objects. Once generation is done, you may re-assemble the dataset into a DataFrame.

Example:

```
my_batch.generate_all_batch_lines()
# Wait for all generation to complete
synthetic_df = my_batch.batches_to_df()
```

Parameters

- **max_invalid** – The number of invalid lines, per batch. If this number is exceeded for any batch, generation will stop.
- **raise_on_failed_batch** – If True, then an exception will be raised if any single batch fails to generate the requested number of lines. If False, then the failed batch will be set to False in the result dictionary from this method.
- **num_lines** – The number of lines to create from each batch. If None then the value from the config template will be used.

Note: Will be overridden / ignored if `seed_fields` is a list. Will be set to the len of the list.

- **seed_fields** – A dictionary that maps field/column names to initial seed values for those columns. This seed will only apply to the first batch that gets trained and generated. Additionally, the fields provided in the mapping **MUST** exist at the front of the first batch.

Note: This param may also be a list of dicts. If this is the case, then `num_lines` will automatically be set to the list length downstream, and a 1:1 ratio will be used for generating valid lines for each prefix.

- **parallelism** – The number of concurrent workers to use. 1 (the default) disables parallelization, while a non-positive value means “number of CPUs + x” (i.e., use 0 for using as many workers as there are CPUs). A floating-point value is interpreted as a fraction of the available CPUs, rounded down.

Returns

A dictionary of batch number to a dictionary that reports the number of valid, invalid lines and bool value that shows if each batch was able to generate the full number of requested lines:

```
{
  0: GenerationSummary(valid_lines=1000, invalid_lines=10, is_
    valid=True),
  1: GenerationSummary(valid_lines=500, invalid_lines=5, is_
    valid=True)
}
```

generate_batch_lines(*batch_idx: int, max_invalid=1000, raise_on_exceed_invalid: bool = False, num_lines: int | None = None, seed_fields: dict | List[dict] | None = None, parallelism: int = 0*) → *GenerationSummary*

Generate lines for a single batch. Lines generated are added to the underlying Batch object for each batch. The lines can be accessed after generation and re-assembled into a DataFrame.

Parameters

- **batch_idx** – The batch number
- **max_invalid** – The max number of invalid lines that can be generated, if this is exceeded, generation will stop
- **raise_on_exceed_invalid** – If true and if the number of lines generated exceeds the `max_invalid` amount, we will re-raise the error thrown by the generation module which will interrupt the running process. Otherwise, we will not raise the caught exception and just return `False` indicating that the batch failed to generate all lines.
- **num_lines** – The number of lines to generate, if `None`, then we use the number from the batch’s config
- **seed_fields** – A dictionary that maps field/column names to initial seed values for those columns. This seed will only apply to the first batch that gets trained and generated. Additionally, the fields provided in the mapping **MUST** exist at the front of the first batch.

Note: This param may also be a list of dicts. If this is the case, then `num_lines` will automatically be set to the list length downstream, and a 1:1 ratio will be used for generating valid lines for each prefix.

- **parallelism** – The number of concurrent workers to use. 1 (the default) disables parallelization, while a non-positive value means “number of CPUs + x” (i.e., use 0 for using as many workers as there are CPUs). A floating-point value is interpreted as a fraction of the available CPUs, rounded down.

master_header_list: List[str]

During training, this is the original column order. When reading from disk, we concatenate all headers from all batches together. This list is not guaranteed to preserve the original header order.

original_headers: List[str]

Stores the original header list / order from the original training data that was used. This is written out to the model directory during training and loaded back in when using read-only mode.

set_batch_validator(*batch_idx: int, validator: Callable*)

Set a validator for a specific batch. If a validator is configured for a batch, each generated record from that batch will be sent to the validator.

Parameters

- **batch_idx** – The batch number .
- **validator** – A callable that should take exactly one argument, which will be the raw line generated from the `generate_text` function.

train_all_batches()

Train a model for each batch.

train_batch(*batch_idx: int*)

Train a model for a single batch. All model information will be written into that batch’s directory.

Parameters

- **batch_idx** – The index of the batch, from the `batches` dictionary

```
class gretel_synthetics.batch.GenerationProgress(current_valid_count: int = 0, current_invalid_count:
                                                int = 0, new_valid_count: int = 0,
                                                new_invalid_count: int = 0, completion_percent:
                                                float = 0.0, timestamp: float = <factory>)
```

This class should not have to be used directly.

It is used to communicate the current progress of record generation.

When a callback function is passed to the `RecordFactory.generate_all()` method, each time the callback is called an instance of this class will be passed as the single argument:

```
def my_callback(data: GenerationProgress):
    ...

factory: RecordFactory
df = factory.generate_all(output="df", callback=my_callback)
```

This class is used to periodically communicate progress of generation to the user, through a callback that can be passed to `RecordFactory.generate_all()` method.

completion_percent: float = 0.0

The percentage of valid lines/records that have been generated.

current_invalid_count: int = 0

The number of invalid lines/records that were generated so far.

current_valid_count: `int = 0`

The number of valid lines/records that were generated so far.

new_invalid_count: `int = 0`

The number of new valid lines/records that were generated since the last progress callback.

new_valid_count: `int = 0`

The number of new valid lines/records that were generated since the last progress callback.

timestamp: `float`

The timestamp from when the information in this object has been captured.

class `gretel_synthetics.batch.GenerationResult`(*records: pandas.core.frame.DataFrame | List[dict],*
exception: Exception | None = None)

class `gretel_synthetics.batch.GenerationSummary`(*valid_lines: int = 0, invalid_lines: int = 0, is_valid:*
bool = False)

A class to capture the summary data after synthetic data is generated.

class `gretel_synthetics.batch.RecordFactory`(**, num_lines: int, batches: dict, header_list: list, delimiter:*
str, seed_fields: dict | list | None = None,
max_invalid=1000, validator: Callable | None = None,
parallelism: int = 4, invalid_cache_size: int = 100)

A stateful factory that can be used to generate and validate entire records, regardless of the number of underlying header clusters that were used to build multiple sub-models.

Instances of this class should be created by calling the appropriate method of the `DataFrameBatch` instance. This class should not have to be used directly. You should be able to create an instance like so:

```
factory = batcher.create_record_factory(num_lines=50)
```

The class is init'd with default capacity and limits as specified by the `num_lines` and `max_invalid` attributes. At any time, you can inspect the state of the instance by doing:

```
factory.summary
```

The factory instance can be used one of two ways: buffered or unbuffered.

For unbuffered mode, the entire instance can be used as an iterator to create synthetic records. Each record will be a dictionary.

Note: All values in the generated dictionaries will be strings.

The `valid_count` and `invalid_count` counters will update as records are generated.

When creating the record factory, you may also provide an entire record validator:

```
def validator(rec: dict):  
    ...  
  
factory = batcher.create_record_factory(num_lines=50, validator=validator)
```

Each generated record dict will be passed to the validator. This validator may either return `False` or raise an exception to mark a record as invalid.

At any point, you may reset the state of the factory by calling:

```
factory.reset()
```

This will reset all counters and allow you to keep generating records.

Finally, you can generate records in buffered mode, where generated records will be buffered in memory and returned as one collection. By default, a list of dicts will be returned:

```
factory.generate_all()
```

You may request the records to be returned as a DataFrame. The dtypes will be inferred as if you were reading the data from a CSV:

```
factory.generate_all(output="df")
```

Note: When using `generate_all`, the factory states will be reset automatically.

generate_all(*output: str | None = None, callback: callable | None = None, callback_interval: int = 30, callback_threading: bool = False*) → *GenerationResult*

Attempt to generate the full number of records that was set when creating the `RecordFactory`. This method will create a buffer that holds all records and then returns the the buffer once generation is complete.

Parameters

- **output** – How the records should be returned. If `None`, which is the default, then a list of record dicts will be returned. Other options that are supported are: 'df' for a DataFrame.
- **callback** – An optional callable that will periodically be called with a `GenerationProgress` instance as the single argument while records are being generated.
- **callback_interval** – If using a callback, the minimum number of seconds that should occur between callbacks.
- **callback_threading** – If enabled, a watchdog thread will be used to execute the callback. This will ensure that the callback is called regardless of invalid or valid counts. If callback threading is disabled, the callback will only be called after valid records are generated. If the callback raises an exception, then a threading event will be set which will trigger the stopping of generation.

Returns

Generated records in an object that is dependent on the `output` param. By default this will be a list of dicts.

validator: Callable

An optional callable that will receive a fully constructed record for one final validation before returning or yielding a single record. Records that do not pass this validation will also increment the `invalid_count`.

7.6 Utils

The utils module provides a number of different methods that are useful for training and working with synthetic data.

Some of these methods carry heavy dependencies such as scikit-learn. To prevent adding unnecessary requirements to the main gretel-synthetics package, util dependencies are shipped under an extra called, `utils`. To install the `utils` extra, you may run

```
pip install -U gretel-synthetics[utils]
```

7.6.1 Stats

Generates correlation reports between data sets.

`gretel_synthetics.utils.stats.calculate_correlation(df: DataFrame, nominal_columns: List[str] | None = None, job_count: int = 4, opt: bool = False) → DataFrame`

Given a dataframe, calculate a matrix of the correlations between the various rows. We use the `calculate_pearsons_r`, `calculate_correlation_ratio` and `calculate_theils_u` to fill in the matrix values.

Parameters

- **df** – The input dataframe.
- **nominal_columns** – Columns to treat as categorical.
- **job_count** – For parallelization of computations.
- **opt** – “optimized.” If `opt` is `True`, then go the faster (just not quite as accurate) route of global replace missing with 0.

Returns

A dataframe of correlation values.

`gretel_synthetics.utils.stats.calculate_correlation_ratio(x, y, opt)`

Calculates the Correlation Ratio for categorical-continuous association. Used in constructing correlation matrix. See http://shakedzy.xyz/dython/modules/nominal/#correlation_ratio.

Parameters

- **x** – first input array, categorical.
- **y** – second input array, numeric.
- **opt** – “optimized.” If `False`, drop missing values if `y` (the numeric column) is `null/nan`.

Returns

float in the range of `[0,1]`.

`gretel_synthetics.utils.stats.calculate_pearsons_r(x, y, opt) → Tuple[float, float]`

Calculate the Pearson correlation coefficient for this pair of rows of our correlation matrix. See <https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.pearsonr.html>.

Parameters

- **x** – first input array.
- **y** – second input array.
- **opt** – “optimized.” If `False`, drop missing values when either the `x` or `y` value is `null/nan`. If `True`, we’ve already replaced `nan`’s with 0’s for entire datafile.

Returns

As per scipy, tuple of Pearson's correlation coefficient and Two-tailed p-value.

`gretel_synthetics.utils.stats.calculate_theils_u(x, y)`

Calculates Theil's U statistic (Uncertainty coefficient) for categorical-categorical association. Used in constructing correlation matrix. See http://shakedzy.xyz/dython/modules/nominal/#theils_u.

Parameters

- **x** – first input array, categorical.
- **y** – second input array, categorical.

Returns

float in the range of [0,1].

`gretel_synthetics.utils.stats.compute_distribution_distance(d1: dict, d2: dict) → float`

Calculates the Jensen Shannon distance between two distributions.

Parameters

- **d1** – Distribution dict. Values must be a probability vector (all values are floats in [0,1], sum of all values is 1.0).
- **d2** – Another distribution dict.

Returns

The distance between the two vectors, range in [0, 1].

Return type

float

`gretel_synthetics.utils.stats.compute_pca(df: DataFrame, n_components: int = 2) → DataFrame`

Do PCA on a dataframe. See <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>.

Parameters

- **df** – The dataframe to analyze for principal components.
- **n_components** – Number of components to keep.

Returns

Dataframe of principal components.

`gretel_synthetics.utils.stats.count_memorized_lines(df1: DataFrame, df2: DataFrame) → int`

Checks for overlap between training and synthesized data.

Parameters

- **df1** – DataFrame of training data.
- **df2** – DataFrame of synthetic data.

Returns

int, the number of overlapping elements.

`gretel_synthetics.utils.stats.get_categorical_field_distribution(field: Series) → dict`

Calculates the normalized distribution of a categorical field.

Parameters

field – A sanitized column extracted from one of the df's.

Returns

keys are the unique values in the field, values are percentages (floats in [0, 100]).

Return type

dict

`gretel_synthetics.utils.stats.get_numeric_distribution_bins(training: Series, synthetic: Series)`

To calculate the distribution distance between two numeric series a la categorical fields we need to bin the data. We want the same bins between both series, based on scrubbed data.

Parameters

- **training** – The numeric series from the training dataframe.
- **synthetic** – The numeric series from the synthetic dataframe.

Returns

bin_edges, numpy array of dtype float

`gretel_synthetics.utils.stats.get_numeric_field_distribution(field: Series, bins) → dict`

Calculates the normalized distribution of a numeric field cut into bins.

Parameters

- **field** – A sanitized column extracted from one of the df's.
- **bins** – Usually an np.ndarray from get_bins, but can be anything that can be safely passed to pandas.cut.

Returns

keys are the unique values in the field, values are floats in [0, 1].

Return type

dict

`gretel_synthetics.utils.stats.normalize_dataset(df: DataFrame) → DataFrame`

Prep a dataframe for PCA. Divide the dataframe into numeric and categorical, fill missing values and encode categorical columns by the frequency of each value and standardize all values.

Parameters**df** – The dataframe to be subjected to PCA.**Returns**

The dataframe, normalized.

7.6.2 Header Clusters

```
gretel_synthetics.utils.header_clusters.cluster(df: DataFrame, header_prefix: List[str] | None =
None, maxsize: int = 20,
average_record_length_threshold: float = 0, method:
str = 'single', numeric_cat: List[str] | None = None,
plot: bool = False, isolate_complex_field: bool =
True) → List[List[str]]
```

Given an input dataframe, extract clusters of similar headers based on a set of heuristics. :param df: The dataframe to cluster headers from. :param header_prefix: List of columns to remove before cluster generation. :param maxsize: The max number of fields in a cluster. :param average_record_length_threshold: Threshold for how long a cluster's records can be.

The default, 0, turns off the average record length (arl) logic. To use arl, use a positive value. Based on our research we recommend setting this value to 250.0.

Parameters

- **method** – Linkage method used to compute header cluster distances. For more information please refer to the scipy docs, <https://docs.scipy.org/doc/scipy/reference/generated/scipy.cluster.hierarchy.linkage.html#scipy-cluster-hierarchy-linkage>. # noqa
- **numeric_cat** – A list of fields to define as categorical. The header clustering code will automatically define pandas “object” and “category” columns as categorical. The `numeric_cat` parameter may be used to define additional categorical fields that may not automatically get identified as such.
- **plot** – Plot header list as a dendrogram.
- **isolate_complex_field** – Enables isolation of complex fields when clustering.

Returns

A list of lists of column names, each column name list being an identified cluster.

7.7 Timeseries DGAN

The Timeseries DGAN module contains a PyTorch implementation of the DoppelGANger model, see <https://arxiv.org/abs/1909.13403> for a detailed description of the model.

```
import numpy as np
from gretel_synthetics.timeseries_dgan.dgan import DGAN
from gretel_synthetics.timeseries_dgan.config import DGANConfig

attributes = np.random.rand(10000, 3)
features = np.random.rand(10000, 20, 2)

config = DGANConfig(
    max_sequence_len=20,
    sample_len=5,
    batch_size=1000,
    epochs=10
)
model = DGAN(config)

model.train(attributes, features)

synthetic_attributes, synthetic_features = model.generate(1000)
```

```
class gretel_synthetics.timeseries_dgan.config.DGANConfig(max_sequence_len: int, sample_len: int,
    attribute_noise_dim: int = 10,
    feature_noise_dim: int = 10,
    attribute_num_layers: int = 3,
    attribute_num_units: int = 100,
    feature_num_layers: int = 1,
    feature_num_units: int = 100,
    use_attribute_discriminator: bool =
    True, normalization: Normalization =
    Normalization.ZERO_ONE,
    apply_feature_scaling: bool = True,
    apply_example_scaling: bool = True,
    binary_encoder_cutoff: int = 150,
    forget_bias: bool = False,
    gradient_penalty_coef: float = 10.0,
    attribute_gradient_penalty_coef: float =
    10.0, attribute_loss_coef: float = 1.0,
    generator_learning_rate: float = 0.001,
    generator_beta1: float = 0.5,
    discriminator_learning_rate: float =
    0.001, discriminator_beta1: float = 0.5,
    attribute_discriminator_learning_rate:
    float = 0.001,
    attribute_discriminator_beta1: float =
    0.5, batch_size: int = 1024, epochs: int
    = 400, discriminator_rounds: int = 1,
    generator_rounds: int = 1, cuda: bool =
    True, mixed_precision_training: bool =
    False)
```

Config object with parameters for training a DGAN model.

Parameters

- **max_sequence_len** – length of time series sequences, variable length sequences are not supported, so all training and generated data will have the same length sequences
- **sample_len** – time series steps to generate from each LSTM cell in DGAN, must be a divisor of max_sequence_len
- **attribute_noise_dim** – length of the GAN noise vectors for attribute generation
- **feature_noise_dim** – length of GAN noise vectors for feature generation
- **attribute_num_layers** – # of layers in the GAN discriminator network
- **attribute_num_units** – # of units per layer in the GAN discriminator network
- **feature_num_layers** – # of LSTM layers in the GAN generator network
- **feature_num_units** – # of units per layer in the GAN generator network
- **use_attribute_discriminator** – use separate discriminator only on attributes, helps DGAN match attribute distributions, Default: True
- **normalization** – default normalization for continuous variables, used when metadata output is not specified during DGAN initialization
- **apply_feature_scaling** – scale each continuous variable to [0,1] or [-1,1] (based on normalization param) before training and rescale to original range during generation, if False

then training data must be within range and DGAN will only generate values in [0,1] or [-1,1], Default: True

- **apply_example_scaling** – compute midpoint and halfrange (equivalent to min/max) for each time series variable and include these as additional attributes that are generated, this provides better support for time series with highly variable ranges, e.g., in network data, a dial-up connection has bandwidth usage in [1kb, 10kb], while a fiber connection is in [100mb, 1gb], Default: True
- **binary_encoder_cutoff** – use binary encoder (instead of one hot encoder) for any column with more than this many unique values. This helps reduce memory consumption for datasets with a lot of unique values.
- **forget_bias** – initialize forget gate bias paramters to 1 in LSTM layers, when True initialization matches tf1 LSTMCell behavior, otherwise default pytorch initialization is used, Default: False
- **gradient_penalty_coef** – coefficient for gradient penalty in Wasserstein loss, Default: 10.0
- **attribute_gradient_penalty_coef** – coefficient for gradient penalty in Wasserstein loss for the attribute discriminator, Default: 10.0
- **attribute_loss_coef** – coefficient for attribute discriminator loss in comparison the standard discriminator on attributes and features, higher values should encourage DGAN to learn attribute distributions, Default: 1.0
- **generator_learning_rate** – learning rate for Adam optimizer
- **generator_beta1** – Adam param for exponential decay of 1st moment
- **discriminator_learning_rate** – learning rate for Adam optimizer
- **discriminator_beta1** – Adam param for exponential decay of 1st moment
- **attribute_discriminator_learning_rate** – learning rate for Adam optimizer
- **attribute_discriminator_beta1** – Adam param for exponential decay of 1st moment
- **batch_size** – # of examples used in batches, for both training and generation
- **epochs** – # of epochs to train model discriminator_rounds: training steps
- **discriminator** (*for the*) – batch
- **generator_rounds** – training steps for the generator in each batch
- **cuda** – use GPU if available
- **mixed_precision_training** – enabling automatic mixed precision while training in order to reduce memory costs, bandwith, and time by identifying the steps that require full precision and using 32-bit floating point for only those steps while using 16-bit floating point everywhere else.

to_dict()

Return dictionary representation of DGANConfig.

Returns

Dictionary of member variables, usable to initialize a new config object, e.g., *DGANConfig(**config.to_dict())*

class gretel_synthetics.timeseries_dgan.config.DfStyle(value)

Supported styles for parsing pandas DataFrames.

See `train_dataframe` method in `dgan.py` for details.

class gretel_synthetics.timeseries_dgan.config.Normalization(value)

Normalization types for continuous variables.

Determines if a sigmoid (ZERO_ONE) or tanh (MINUSONE_ONE) activation is used for the output layers in the generation network.

class gretel_synthetics.timeseries_dgan.config.OutputType(value)

Supported variables types.

Determines internal representation of variables and output layers in generation network.

PyTorch implementation of DoppelGANger, from <https://arxiv.org/abs/1909.13403>

Based on tensorflow 1 code in <https://github.com/fjxmlzn/DoppelGANger>

DoppelGANger is a generative adversarial network (GAN) model for time series. It supports multi-variate time series (referred to as features) and fixed variables for each time series (attributes). The combination of attribute values and sequence of feature values is 1 example. Once trained, the model can generate novel examples that exhibit the same temporal correlations as seen in the training data. See <https://arxiv.org/abs/1909.13403> for additional details on the model.

As a reference for terminology, consider open-high-low-close (OHLC) data from stock markets. Each stock is an example, with fixed attributes such as exchange, sector, country. The features or time series consists of open, high, low, and closing prices for each time interval (daily). After being trained on historical data, the model can generate more hypothetical stocks and price behavior on the training time range.

Sample usage:

```
import numpy as np
from gretel_synthetics.timeseries_dgan.dgan import DGAN
from gretel_synthetics.timeseries_dgan.config import DGANConfig

attributes = np.random.rand(10000, 3)
features = np.random.rand(10000, 20, 2)

config = DGANConfig(
    max_sequence_len=20,
    sample_len=5,
    batch_size=1000,
    epochs=10
)

model = DGAN(config)

model.train_numpy(attributes=attributes, features=features)

synthetic_attributes, synthetic_features = model.generate_numpy(1000)
```

class gretel_synthetics.timeseries_dgan.dgan.DGAN(config: DGANConfig, attribute_outputs: List[Output] | None = None, feature_outputs: List[Output] | None = None)

DoppelGANger model.

Interface for training model and generating data based on configuration in an DGANConfig instance.

DoppelGANger uses a specific internal representation for data which is hidden from the user in the public interface. Standard usage of DGAN instances should pass continuous variables as floats in the original space (not normalized), and discrete variables may be strings, integers, or floats. This is the format expected by both `train_numpy()` and `train_dataframe()` and the `generate_numpy()` and `generate_dataframe()` functions will return data in this same format. In standard usage, the detailed transformation info in `attribute_outputs` and `feature_outputs` are not needed, those will be created automatically when a `train*` function is called with data.

If more control is needed and you want to use the normalized values and one-hot encoding directly, use the `_train()` and `_generate()` functions. `transformations.py` contains internal helper functions for working with the Output metadata instances and converting data to and from the internal representation. To dive even deeper into the model structure, see the `torch_modules.py` which contains the torch implementations of the networks used in DGAN. As internal details, `transformations.py` and `torch_modules.py` are not part of the public interface and may change at any time without notice.

__init__ (*config: [DGANConfig](#), attribute_outputs: List[Output] | None = None, feature_outputs: List[Output] | None = None*)

Create a DoppelGANger model.

Parameters

- **config** – [DGANConfig](#) containing model parameters
- **attribute_outputs** – custom metadata for attributes, not needed for standard usage
- **feature_outputs** – custom metadata for features, not needed for standard usage

generate_dataframe (*n: int | None = None, attribute_noise: Tensor | None = None, feature_noise: Tensor | None = None*) → DataFrame

Generate synthetic data from DGAN model.

Once trained, a DGAN model can generate arbitrary amounts of synthetic data by sampling from the noise distributions. Specify either the number of records to generate, or the specific noise vectors to use.

Parameters

- **n** – number of examples to generate
- **attribute_noise** – noise vectors to create synthetic data
- **feature_noise** – noise vectors to create synthetic data

Returns

pandas DataFrame in same format used in ‘train_dataframe’ call

generate_numpy (*n: int | None = None, attribute_noise: Tensor | None = None, feature_noise: Tensor | None = None*) → Tuple[ndarray | None, ndarray]

Generate synthetic data from DGAN model.

Once trained, a DGAN model can generate arbitrary amounts of synthetic data by sampling from the noise distributions. Specify either the number of records to generate, or the specific noise vectors to use.

Parameters

- **n** – number of examples to generate
- **attribute_noise** – noise vectors to create synthetic data
- **feature_noise** – noise vectors to create synthetic data

Returns

Tuple of attributes and features as numpy arrays.

classmethod `load(file_name: str, **kwargs) → DGAN`

Load DGAN model instance from a file.

Parameters

- **file_name** – location to load from
- **kwargs** – additional parameters passed to `torch.load`, for example, use `map_location=torch.device("cpu")` to load a model saved for GPU on a machine without cuda

Returns

DGAN model instance

save(`file_name: str, **kwargs`)

Save DGAN model to a file.

Parameters

- **file_name** – location to save serialized model
- **kwargs** – additional parameters passed to `torch.save`

train_dataframe(`df: DataFrame, attribute_columns: List[str] | None = None, feature_columns: List[str] | None = None, example_id_column: str | None = None, time_column: str | None = None, discrete_columns: List[str] | None = None, df_style: DfStyle = DfStyle.WIDE, progress_callback: Callable[[ProgressInfo], None] | None = None`) → None

Train DGAN model on data in pandas DataFrame.

Training data can be in either “wide” or “long” format. “Wide” format uses one row for each example with 0 or more attribute columns and 1 column per time point in the time series. “Wide” format is restricted to 1 feature variable. “Long” format uses one row per time point, supports multiple feature variables, and uses additional example id to split into examples and time column to sort.

Parameters

- **df** – DataFrame of training data
- **attribute_columns** – list of column names containing attributes, if None, no attribute columns are used. Must be disjoint from the feature columns.
- **feature_columns** – list of column names containing features, if None all non-attribute columns are used. Must be disjoint from attribute columns.
- **example_id_column** – column name used to split “long” format data frame into multiple examples, if None, data is treated as a single example. This value must be unique from the other column list parameters.
- **time_column** – column name used to sort “long” format data frame, if None, data frame order of rows/time points is used. This value must be unique from the other column list parameters.
- **discrete_columns** – column names (either attributes or features) to treat as discrete (use one-hot or binary encoding), any string or object columns are automatically treated as discrete
- **df_style** – str enum of “wide” or “long” indicating format of the DataFrame

train_numpy(`features: ndarray, feature_types: List[OutputType] | None = None, attributes: ndarray | None = None, attribute_types: List[OutputType] | None = None, progress_callback: Callable[[ProgressInfo], None] | None = None`) → None

Train DGAN model on data in numpy arrays.

Training data is passed in 2 numpy arrays, one for attributes (2d) and one for features (3d). This data should be in the original space and is not transformed. If the data is already transformed into the internal DGAN representation (continuous variable scaled to [0,1] or [-1,1] and discrete variables one-hot or binary encoded), use the internal `_train()` function instead of `train_numpy()`.

In standard usage, `attribute_types` and `feature_types` may be provided on the first call to `train()` to setup the model structure. If not specified, the default is to assume continuous variables for floats and integers, and discrete for strings. If outputs metadata was specified when the instance was initialized or `train()` was previously called, then `attribute_types` and `feature_types` are not needed.

Parameters

- **features** – 3-d numpy array of time series features for the training, size is (# of training examples) X `max_sequence_len` X (# of features)
- **feature_types** (*Optional*) – Specification of Discrete or Continuous type for each variable of the features. If `None`, assume continuous variables for floats and integers, and discrete for strings. Ignored if the model was already built, either by passing *output params at initialization or because train_* was called previously.
- **attributes** (*Optional*) – 2-d numpy array of attributes for the training examples, size is (# of training examples) X (# of attributes)
- **attribute_types** (*Optional*) – Specification of Discrete or Continuous type for each variable of the attributes. If `None`, assume continuous variables for floats and integers, and discrete for strings. Ignored if the model was already built, either by passing *output params at initialization or because train_* was called previously.

`gretel_synthetics.timeseries_dgan.dgan.find_max_consecutive_nans(array: array) → int`

Returns the maximum number of consecutive NaNs in an array.

Parameters

array – 1-d numpy array of time series per example.

Returns

The maximum number of consecutive NaNs in a times series array.

Return type

`max_cons_nan`

`gretel_synthetics.timeseries_dgan.dgan.nan_linear_interpolation(arrays: ndarray) → ndarray`

Replaces all NaNs via linear interpolation.

Parameters

- **arrays** – 3-d numpy array of continuous features, with shape
- (**#examples** –
- **max_sequence_length** –
- **features**) (*#continuous*) –

Returns

3-d numpy array where NaNs are replaced via linear interpolation.

Return type

`arrays`

```
gretel_synthetics.timeseries_dgan.dgan.validation_check(array: ndarray,  
                                                         invalid_examples_ratio_cutoff: float = 0.5,  
                                                         nans_ratio_cutoff: float = 0.1,  
                                                         consecutive_nans_max: int = 5,  
                                                         consecutive_nans_ratio_cutoff: float =  
                                                         0.05) → array
```

Checks if continuous features of examples are valid.

Returns a 1-d numpy array of booleans with shape (#examples) indicating valid examples. Examples with continuous features fall into 3 categories: good, valid (fixable) and invalid (non-fixable). - “Good” examples have no NaNs. - “Valid” examples have a low percentage of nans and a below a threshold number of consecutive NaNs. - “Invalid” are the rest, and are marked “False” in the returned array. Later on, these are omitted from training. If there are too many, later, we error out.

Parameters

- **array** – 3-d numpy array of continuous features with
- **shape** (*#examples, max_sequence_length, #continuous features*) –
- **invalid_examples_ratio_cutoff** – Error out if the invalid examples ratio in the dataset
- **value.** (*is higher than this*) –
- **nans_ratio_cutoff** – If the percentage of nans for any continuous feature in an example
- **value** (*is greater than this*) –
- **invalid.** (*then the example is*) –
- **consecutive_nans_max** – If the maximum number of consecutive nans in a continuous
- **number** (*feature is greater than this*) –
- **invalid.** –
- **consecutive_nans_ratio_cutoff** – If the maximum number of consecutive nans in a
- **example** (*continuous feature is greater than this ratio times the length of the*) –
- **samples** (*(number)*) –
- **invalid.** –

Returns

1-d numpy array of booleans indicating valid examples with shape (#examples).

Return type

valid_examples

7.8 ACTGAN

The ACTGAN sub-package contains an alternate implementation of the SDV CTGAN model. It provides some improvement and automation around automatic detection of datetime fields and optional usage of a binary encoder for discrete columns for better memory usage.

Please see the “ACTGAN_Demo” Notebook in the “examples” directory in the repository root.

Wrapper around ACTGAN model.

```

class gretel_synthetics.actgan.actgan_wrapper.ACTGAN(field_names: List[str] | None = None,
                                                    field_types: Dict[str, dict] | None = None,
                                                    field_transformers: Dict[str, BaseTransformer |
str] | None = None, auto_transform_datetimes:
bool = False, anonymize_fields: Dict[str, str] |
None = None, primary_key: str | None = None,
constraints: List[Constraint] | List[dict] | None
= None, table_metadata: Metadata | dict |
None = None, embedding_dim: int = 128,
generator_dim: Sequence[int] = (256, 256),
discriminator_dim: Sequence[int] = (256,
256), generator_lr: float = 0.0002,
generator_decay: float = 1e-06,
discriminator_lr: float = 0.0002,
discriminator_decay: float = 1e-06, batch_size:
int = 500, discriminator_steps: int = 1,
binary_encoder_cutoff: int = 500,
binary_encoder_nan_handler: str | None =
None, cbn_sample_size: int | None = 250000,
log_frequency: bool = True, verbose: bool =
False, epochs: int = 300, epoch_callback:
Callable[[EpochInfo], None] | None = None,
pac: int = 10, cuda: bool = True,
learn_rounding_scheme: bool = True,
enforce_min_max_values: bool = True,
conditional_vector_type:
ConditionalVectorType =
ConditionalVectorType.SINGLE_DISCRETE,
conditional_select_mean_columns: float | None
= None, conditional_select_column_prob: float
| None = None, reconstruction_loss_coef: float
= 1.0, force_conditioning: bool = False)

```

Parameters

- **field_names** – List of names of the fields that need to be modeled and included in the generated output data. Any additional fields found in the data will be ignored and will not be included in the generated output. If `None`, all the fields found in the data are used.
- **field_types** – Dictionary specifying the data types and subtypes of the fields that will be modeled. Field types and subtypes combinations must be compatible with the SDV Metadata Schema.
- **field_transformers** – Dictionary specifying which transformers to use for each field. Available transformers are:
 - `FloatFormatter`: Uses a `FloatFormatter` for numerical data.
 - `FrequencyEncoder`: Uses a `FrequencyEncoder` without gaussian noise.
 - `FrequencyEncoder_noised`: Uses a `FrequencyEncoder` adding gaussian noise.
 - `OneHotEncoder`: Uses a `OneHotEncoder`.
 - `LabelEncoder`: Uses a `LabelEncoder` without gaussian noise.
 - `LabelEncoder_noised`: Uses a `LabelEncoder` adding gaussian noise.
 - `BinaryEncoder`: Uses a `BinaryEncoder`.

– `UnixTimestampEncoder`: Uses a `UnixTimestampEncoder`.

NOTE: Specifically for ACTGAN, some attributes such as `auto_transform_datetimes` will automatically attempt to detect field types and will automatically set the `field_transformers` dictionary at construction time. However, autodetection of `field_types` and `field_transformers` will not be over-written by any concrete values that were provided to this constructor.

- **`auto_transform_datetimes`** – If set, prior to fitting, each column will be checked for being a potential “datetime” type. For each column that is discovered as a “datetime” the `field_types` and `field_transformers` SDV metadata dicts will be automatically updated such that datetimes are transformed to Unix timestamps. NOTE: if fields are already specified in `field_types` or `field_transformers` these fields will be skipped by the auto detector.
- **`anonymize_fields`** – Dict specifying which fields to anonymize and what faker category they belong to.
- **`primary_key`** – Name of the field which is the primary key of the table.
- **`constraints`** – List of Constraint objects or dicts.
- **`table_metadata`** – Table metadata instance or dict representation. If given alongside any other metadata-related arguments, an exception will be raised. If not given at all, it will be built using the other arguments or learned from the data.
- **`embedding_dim`** – Size of the random sample passed to the Generator. Defaults to 128.
- **`generator_dim`** – Size of the output samples for each one of the Residuals. A Residual Layer will be created for each one of the values provided. Defaults to (256, 256).
- **`discriminator_dim`** – Size of the output samples for each one of the Discriminator Layers. A Linear Layer will be created for each one of the values provided. Defaults to (256, 256).
- **`generator_lr`** – Learning rate for the generator. Defaults to 2e-4.
- **`generator_decay`** – Generator weight decay for the Adam Optimizer. Defaults to 1e-6.
- **`discriminator_lr`** – Learning rate for the discriminator. Defaults to 2e-4.
- **`discriminator_decay`** – Discriminator weight decay for the Adam Optimizer. Defaults to 1e-6.
- **`batch_size`** – Number of data samples to process in each step.
- **`discriminator_steps`** – Number of discriminator updates to do for each generator update. From the WGAN paper: <https://arxiv.org/abs/1701.07875>. WGAN paper default is 5. Default used is 1 to match original CTGAN implementation.
- **`binary_encoder_cutoff`** – For any given column, the number of unique values that should exist before switching over to binary encoding instead of OHE. This will help reduce memory consumption for datasets with a lot of unique values.
- **`binary_encoder_nan_handler`** – Binary encoding currently may produce errant NaN values during reverse transformation. By default these NaN’s will be left in place, however if this value is set to “mode” then those NaN’ will be replaced by a random value that is a known mode for a given column.
- **`cbn_sample_size`** – Number of rows to sample from each column for identifying clusters for the cluster-based normalizer. This only applies to float columns. If set to 0, no sampling is done and all values are considered, which may be very slow. Defaults to 250_000.
- **`log_frequency`** – Whether to use log frequency of categorical levels in conditional sampling. Defaults to True.

- **verbose** – Whether to have print statements for progress results. Defaults to False.
- **epochs** – Number of training epochs. Defaults to 300.
- **epoch_callback** – An optional function to call after each epoch, the argument will be a `EpochInfo` instance
- **pac** – Number of samples to group together when applying the discriminator. Defaults to 10.
- **cuda** – If True, use CUDA. If a str, use the indicated device. If False, do not use cuda at all. Defaults to True.
- **learn_rounding_scheme** – Define rounding scheme for `FloatFormatter`. If True, the data returned by `reverse_transform` will be rounded to that place. Defaults to True.
- **enforce_min_max_values** – Specify whether or not to clip the data returned by `reverse_transform` of the numerical transformer, `FloatFormatter`, to the min and max values seen during `fit`. Defaults to True.
- **conditional_vector_type** – Type of conditional vector to include in input to the generator. Influences how effective and flexible the native conditional generation is. Options include `SINGLE_DISCRETE` (original CTGAN setup) and `ANYWAY`. Default is `SINGLE_DISCRETE`.
- **conditional_select_mean_columns** – Target number of columns to select for conditioning on average during training. Only used for `ANYWAY` conditioning. Use if typical number of columns to seed on is known. If set, `conditional_select_column_prob` must be None. Equivalent to setting `conditional_select_column_prob` to `conditional_select_mean_columns / # of columns`. Defaults to None.
- **conditional_select_column_prob** – Probability to select any given column to be conditioned on during training. Only used for `ANYWAY` conditioning. If set, `conditional_select_mean_columns` must be None. Defaults to None.
- **reconstruction_loss_coef** – Multiplier on reconstruction loss, higher values focus the generator optimization more on accurate conditional vector generation. Defaults to 1.0.
- **force_conditioning** – Directly set the requested conditional generation columns in generated data. Will bypass rejection sampling and be faster, but may reduce quality of the generated data and correlations between conditioned columns and other variables may be weaker. Defaults to False.

fit(*args, **kwargs)

Fit the ACTGAN model to the provided data. Prior to fitting, specific auto-detection of data types will be done if the provided data is a `DataFrame`.

sample(*args, **kwargs)

Sample rows from this table.

Parameters

- **num_rows** (*int*) – Number of rows to sample. This parameter is required.
- **randomize_samples** (*bool*) – Whether or not to use a fixed seed when sampling. Defaults to True.
- **max_tries_per_batch** (*int*) – Number of times to retry sampling until the batch size is met. Defaults to 100.
- **batch_size** (*int* or *None*) – The batch size to sample. Defaults to `num_rows`, if None.

- **output_file_path** (*str or None*) – The file to periodically write sampled rows to. If *None*, does not write rows anywhere.
- **conditions** – Deprecated argument. Use the *sample_conditions* method with *sdv.sampling.Condition* objects instead.

Returns

Sampled data.

Return type

`pandas.DataFrame`

sample_remaining_columns(*args, **kwargs)

Sample rows from this table.

Parameters

- **known_columns** (*pandas.DataFrame*) – A *pandas.DataFrame* with the columns that are already known. The output is a *DataFrame* such that each row in the output is sampled conditionally on the corresponding row in the input.
- **max_tries_per_batch** (*int*) – Number of times to retry sampling until the batch size is met. Defaults to 100.
- **batch_size** (*int*) – The batch size to use per sampling call.
- **randomize_samples** (*bool*) – Whether or not to use a fixed seed when sampling. Defaults to *True*.
- **output_file_path** (*str or None*) – The file to periodically write sampled rows to. Defaults to a temporary file, if *None*.

Returns

Sampled data.

Return type

`pandas.DataFrame`

Raises

- **ConstraintsNotMetError** – If the conditions are not valid for the given constraints.
- **ValueError** – If any of the following happens: * any of the conditions' columns are not valid. * no rows could be generated.

Complex datastructures for ACTGAN

```
class gretel_synthetics.actgan.structures.ColumnIdInfo(discrete_column_id: 'int', column_id: 'int',
                                                    value_id: 'np.ndarray')
```

```
class gretel_synthetics.actgan.structures.ColumnTransformInfo(column_name: 'str', column_type:
                                                            'ColumnType', transform:
                                                            'BaseTransformer', encodings:
                                                            'List[ColumnEncoding]')
```

```
class gretel_synthetics.actgan.structures.ColumnType(value)
    An enumeration.
```

```
class gretel_synthetics.actgan.structures.ConditionalVectorType(value)
    An enumeration.
```

```
class gretel_synthetics.actgan.structures.EpochInfo(epoch: int, loss_g: float, loss_d: float, loss_r:  
                                                    float)
```

When creating a model such as ACTGAN if the `epoch_callback` attribute is set to a callable, then after each epoch the provided callable will be called with an instance of this class as the only argument.

INDICES AND TABLES

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

PYTHON MODULE INDEX

g

- `gretel_synthetics.actgan.actgan_wrapper`, 42
- `gretel_synthetics.actgan.structures`, 46
- `gretel_synthetics.batch`, 25
- `gretel_synthetics.config`, 15
- `gretel_synthetics.generate`, 22
- `gretel_synthetics.timeseries_dgan.config`, 35
- `gretel_synthetics.timeseries_dgan.dgan`, 38
- `gretel_synthetics.tokenizers`, 19
- `gretel_synthetics.train`, 22
- `gretel_synthetics.utils.header_clusters`, 34
- `gretel_synthetics.utils.stats`, 32

Symbols

`__init__()` (*gretel_synthetic.timeseries_dgan.dgan.DGAN* method), 39

A

`ACTGAN` (class in *gretel_synthetic.actgan.actgan_wrapper*), 42

`add_valid_data()` (*gretel_synthetic.batch.Batch* method), 25

`annotate_data()` (*gretel_synthetic.tokenizers.BaseTokenizerTrainer* method), 20

`as_dict()` (*gretel_synthetic.config.BaseConfig* method), 15

`as_dict()` (*gretel_synthetic.generate.gen_text* method), 23

B

`Base` (class in *gretel_synthetic.tokenizers*), 19

`BaseConfig` (class in *gretel_synthetic.config*), 15

`BaseGenerator` (class in *gretel_synthetic.generate*), 22

`BaseTokenizer` (class in *gretel_synthetic.tokenizers*), 19

`BaseTokenizerTrainer` (class in *gretel_synthetic.tokenizers*), 20

`Batch` (class in *gretel_synthetic.batch*), 25

`batch_size` (*gretel_synthetic.batch.DataFrameBatch* attribute), 26

`batch_to_df()` (*gretel_synthetic.batch.DataFrameBatch* method), 26

`batches` (*gretel_synthetic.batch.DataFrameBatch* attribute), 27

`batches_to_df()` (*gretel_synthetic.batch.DataFrameBatch* method), 27

C

`calculate_correlation()` (in module *gretel_synthetic.utils.stats*), 32

`calculate_correlation_ratio()` (in module *gretel_synthetic.utils.stats*), 32

`calculate_pearsons_r()` (in module *gretel_synthetic.utils.stats*), 32

`calculate_theils_u()` (in module *gretel_synthetic.utils.stats*), 33

`character_coverage` (*gretel_synthetic.tokenizers.SentencePieceTokenizerTrainer* attribute), 21

`CharTokenizer` (class in *gretel_synthetic.tokenizers*), 20

`CharTokenizerTrainer` (class in *gretel_synthetic.tokenizers*), 20

`checkpoint_dir` (*gretel_synthetic.config.BaseConfig* attribute), 15

`cluster()` (in module *gretel_synthetic.utils.header_clusters*), 34

`ColumnIdInfo` (class in *gretel_synthetic.actgan.structures*), 46

`ColumnTransformInfo` (class in *gretel_synthetic.actgan.structures*), 46

`ColumnType` (class in *gretel_synthetic.actgan.structures*), 46

`completion_percent` (*gretel_synthetic.batch.GenerationProgress* attribute), 29

`compute_distribution_distance()` (in module *gretel_synthetic.utils.stats*), 33

`compute_pca()` (in module *gretel_synthetic.utils.stats*), 33

`ConditionalVectorType` (class in *gretel_synthetic.actgan.structures*), 46

`config` (*gretel_synthetic.batch.DataFrameBatch* attribute), 27

`config` (*gretel_synthetic.tokenizers.BaseTokenizerTrainer* attribute), 20

`config_from_model_dir()` (in module *gretel_synthetic.config*), 19

`CONFIG_MAP` (in module *gretel_synthetic.config*), 16

`count_memorized_lines()` (in module *gretel_synthetic.utils.stats*), 33

`create_training_data()` (*gretel_synthetic.batch.DataFrameBatch* method), 27

- `current_invalid_count` (*gretel_synthetics.batch.GenerationProgress* attribute), 29
- `current_valid_count` (*gretel_synthetics.batch.GenerationProgress* attribute), 29
- ## D
- `data_iterator()` (*gretel_synthetics.tokenizers.BaseTokenizerTrainer* method), 20
- `DataFrameBatch` (class in *gretel_synthetics.batch*), 26
- `decode_from_ids()` (*gretel_synthetics.tokenizers.BaseTokenizer* method), 19
- `delimiter` (*gretel_synthetics.generate.gen_text* attribute), 23
- `DfStyle` (class in *gretel_synthetics.timeseries_dgan.config*), 37
- `DGAN` (class in *gretel_synthetics.timeseries_dgan.dgan*), 38
- `DGANConfig` (class in *gretel_synthetics.timeseries_dgan.config*), 35
- ## E
- `encode_to_ids()` (*gretel_synthetics.tokenizers.BaseTokenizer* method), 20
- `epoch_callback` (*gretel_synthetics.config.BaseConfig* attribute), 15
- `EpochInfo` (class in *gretel_synthetics.actgan.structures*), 46
- `EpochState` (class in *gretel_synthetics.train*), 22
- `explain` (*gretel_synthetics.generate.gen_text* attribute), 23
- ## F
- `field_delimiter` (*gretel_synthetics.config.BaseConfig* attribute), 15
- `field_delimiter_token` (*gretel_synthetics.config.BaseConfig* attribute), 15
- `find_max_consecutive_nans()` (in module *gretel_synthetics.timeseries_dgan.dgan*), 41
- `fit()` (*gretel_synthetics.actgan.actgan_wrapper.ACTGAN* method), 45
- ## G
- `gen_text` (class in *gretel_synthetics.generate*), 23
- `generate_all()` (*gretel_synthetics.batch.RecordFactory* method), 31
- `generate_all_batch_lines()` (*gretel_synthetics.batch.DataFrameBatch* method), 27
- `generate_batch_lines()` (*gretel_synthetics.batch.DataFrameBatch* method), 28
- `generate_dataframe()` (*gretel_synthetics.timeseries_dgan.dgan.DGAN* method), 39
- `generate_numpy()` (*gretel_synthetics.timeseries_dgan.dgan.DGAN* method), 39
- `generate_text()` (in module *gretel_synthetics.generate*), 24
- `GenerationProgress` (class in *gretel_synthetics.batch*), 29
- `GenerationResult` (class in *gretel_synthetics.batch*), 30
- `GenerationSummary` (class in *gretel_synthetics.batch*), 30
- `GenText` (class in *gretel_synthetics.generate*), 22
- `get_categorical_field_distribution()` (in module *gretel_synthetics.utils.stats*), 33
- `get_generator_class()` (*gretel_synthetics.config.BaseConfig* method), 15
- `get_generator_class()` (*gretel_synthetics.config.TensorFlowConfig* method), 19
- `get_numeric_distribution_bins()` (in module *gretel_synthetics.utils.stats*), 34
- `get_numeric_field_distribution()` (in module *gretel_synthetics.utils.stats*), 34
- `get_training_callable()` (*gretel_synthetics.config.BaseConfig* method), 15
- `get_training_callable()` (*gretel_synthetics.config.TensorFlowConfig* method), 19
- `get_validator()` (*gretel_synthetics.batch.Batch* method), 25
- `gpu_check()` (*gretel_synthetics.config.BaseConfig* method), 16
- `gpu_check()` (*gretel_synthetics.config.TensorFlowConfig* method), 19
- `gretel_synthetics.actgan.actgan_wrapper` module, 42
- `gretel_synthetics.actgan.structures` module, 46
- `gretel_synthetics.batch` module, 25
- `gretel_synthetics.config` module, 15
- `gretel_synthetics.generate` module, 22

gretel_synthetics.timeseries_dgan.config
 module, 35
 gretel_synthetics.timeseries_dgan.dgan
 module, 38
 gretel_synthetics.tokenizers
 module, 19
 gretel_synthetics.train
 module, 22
 gretel_synthetics.utils.header_clusters
 module, 34
 gretel_synthetics.utils.stats
 module, 32

I

input_data_path (gretel_synthetics.config.BaseConfig attribute), 16

L

load() (gretel_synthetics.timeseries_dgan.dgan.DGAN class method), 39
 load() (gretel_synthetics.tokenizers.BaseTokenizer class method), 20
 load() (gretel_synthetics.tokenizers.CharTokenizer class method), 20
 load() (gretel_synthetics.tokenizers.SentencePieceTokenizer class method), 21
 load_validator_from_file() (gretel_synthetics.batch.Batch method), 25
 LocalConfig (in module gretel_synthetics.config), 16

M

master_header_list (gretel_synthetics.batch.DataFrameBatch attribute), 29
 max_line_line (gretel_synthetics.tokenizers.SentencePieceTokenizer class attribute), 21
 max_lines (gretel_synthetics.config.BaseConfig attribute), 16
 max_training_time_seconds (gretel_synthetics.config.BaseConfig attribute), 16
 model_type (gretel_synthetics.config.BaseConfig attribute), 16
 module
 gretel_synthetics.actgan.actgan_wrapper, 42
 gretel_synthetics.actgan.structures, 46
 gretel_synthetics.batch, 25
 gretel_synthetics.config, 15
 gretel_synthetics.generate, 22
 gretel_synthetics.timeseries_dgan.config, 35
 gretel_synthetics.timeseries_dgan.dgan, 38
 gretel_synthetics.tokenizers, 19
 gretel_synthetics.train, 22
 gretel_synthetics.utils.header_clusters, 34
 gretel_synthetics.utils.stats, 32

N

nan_linear_interpolation() (in module gretel_synthetics.timeseries_dgan.dgan), 41
 new_invalid_count (gretel_synthetics.batch.GenerationProgress attribute), 30
 new_valid_count (gretel_synthetics.batch.GenerationProgress attribute), 30
 Normalization (class in gretel_synthetics.timeseries_dgan.config), 38
 normalize_dataset() (in module gretel_synthetics.utils.stats), 34
 num_lines (gretel_synthetics.tokenizers.BaseTokenizerTrainer attribute), 20

O

original_headers (gretel_synthetics.batch.DataFrameBatch attribute), 29
 OutputType (class in gretel_synthetics.timeseries_dgan.config), 38
 overwrite (gretel_synthetics.config.BaseConfig attribute), 16

P

PredString (in module gretel_synthetics.generate), 22
 pretrain_sentence_count (gretel_synthetics.tokenizers.SentencePieceTokenizerTrainer attribute), 21

R

RecordFactory (class in gretel_synthetics.batch), 30
 reset_gen_data() (gretel_synthetics.batch.Batch method), 25

S

sample() (gretel_synthetics.actgan.actgan_wrapper.ACTGAN method), 45
 sample_remaining_columns() (gretel_synthetics.actgan.actgan_wrapper.ACTGAN method), 46
 save() (gretel_synthetics.timeseries_dgan.dgan.DGAN method), 40
 SeedingGenerator (class in gretel_synthetics.generate), 23
 SentencePieceColumnTokenizer (class in gretel_synthetics.tokenizers), 21

SentencePieceColumnTokenizerTrainer (class in *gretel_synthetics.tokenizers*), 21
 SentencePieceTokenizer (class in *gretel_synthetics.tokenizers*), 21
 SentencePieceTokenizerTrainer (class in *gretel_synthetics.tokenizers*), 21
 set_batch_validator() (gretel_synthetics.batch.DataFrameBatch method), 29
 set_validator() (gretel_synthetics.batch.Batch method), 25
 Settings (class in *gretel_synthetics.generate*), 23
 synthetic_df (gretel_synthetics.batch.Batch property), 26

T

TensorFlowConfig (class in *gretel_synthetics.config*), 16
 text (gretel_synthetics.generate.gen_text attribute), 23
 timestamp (gretel_synthetics.batch.GenerationProgress attribute), 30
 to_dict() (gretel_synthetics.timeseries_dgan.config.DGANConfig method), 37
 tokenizer_from_model_dir() (in module *gretel_synthetics.tokenizers*), 21
 TokenizerError, 21
 total_vocab_size (gretel_synthetics.tokenizers.BaseTokenizer property), 20
 total_vocab_size (gretel_synthetics.tokenizers.CharTokenizer property), 20
 total_vocab_size (gretel_synthetics.tokenizers.SentencePieceTokenizer property), 21
 train() (gretel_synthetics.tokenizers.BaseTokenizerTrainer method), 20
 train() (in module *gretel_synthetics.train*), 22
 train_all_batches() (gretel_synthetics.batch.DataFrameBatch method), 29
 train_batch() (gretel_synthetics.batch.DataFrameBatch method), 29
 train_dataframe() (gretel_synthetics.timeseries_dgan.dgan.DGAN method), 40
 train_numpy() (gretel_synthetics.timeseries_dgan.dgan.DGAN method), 40
 train_rnn() (in module *gretel_synthetics.train*), 22
 training_data_path (gretel_synthetics.config.BaseConfig attribute), 16
 TrainingParams (class in *gretel_synthetics.train*), 22

V

valid (gretel_synthetics.generate.gen_text attribute), 23
 validation_check() (in module *gretel_synthetics.timeseries_dgan.dgan*), 41
 validation_split (gretel_synthetics.config.BaseConfig attribute), 16
 validator (gretel_synthetics.batch.RecordFactory attribute), 31
 values_as_list() (gretel_synthetics.generate.gen_text method), 24
 vocab_size (gretel_synthetics.tokenizers.BaseTokenizerTrainer attribute), 20
 vocab_size (gretel_synthetics.tokenizers.SentencePieceTokenizerTrainer attribute), 21
 VocabSizeTooSmall, 21