
nlpaug Documentation

Release 0.0.16

Edward Ma

Aug 12, 2020

Contents:

| | | |
|----------|----------------------------------------------------------------|-----------|
| 1 | Overview | 3 |
| 2 | Example | 5 |
| 3 | Augmenter | 7 |
| 3.1 | Audio Augmenter | 7 |
| 3.1.1 | nlpaug.augmenter.audio.crop | 7 |
| 3.1.2 | nlpaug.augmenter.audio.loudness | 8 |
| 3.1.3 | nlpaug.augmenter.audio.mask | 9 |
| 3.1.4 | nlpaug.augmenter.audio.noise | 10 |
| 3.1.5 | nlpaug.augmenter.audio.pitch | 11 |
| 3.1.6 | nlpaug.augmenter.audio.shift | 12 |
| 3.1.7 | nlpaug.augmenter.audio.speed | 13 |
| 3.1.8 | nlpaug.augmenter.audio.vt1p | 14 |
| 3.2 | Character Augmenter | 15 |
| 3.2.1 | nlpaug.augmenter.char.keyboard | 15 |
| 3.2.2 | nlpaug.augmenter.char.ocr | 16 |
| 3.2.3 | nlpaug.augmenter.char.random | 18 |
| 3.3 | Sentence Augmenter | 20 |
| 3.3.1 | nlpaug.augmenter.sentence.context_word_embs_sentence | 20 |
| 3.4 | Spectrogram Augmenter | 21 |
| 3.4.1 | nlpaug.augmenter.spectrogram.frequency_masking | 21 |
| 3.4.2 | nlpaug.augmenter.spectrogram.time_masking | 22 |
| 3.5 | Word Augmenter | 23 |
| 3.5.1 | nlpaug.augmenter.word.antonym | 23 |
| 3.5.2 | nlpaug.augmenter.word.context_word_embs | 24 |
| 3.5.3 | nlpaug.augmenter.word.random | 26 |
| 3.5.4 | nlpaug.augmenter.word.spelling | 28 |
| 3.5.5 | nlpaug.augmenter.word.split | 29 |
| 3.5.6 | nlpaug.augmenter.word.synonym | 30 |
| 3.5.7 | nlpaug.augmenter.word.tfidf | 31 |
| 3.5.8 | nlpaug.augmenter.word.word_embs | 33 |
| 4 | Flow | 35 |
| 4.1 | nlpaug.flow.sequential | 35 |
| 4.2 | nlpaug.flow.sometimes | 35 |

| | | |
|----------|-------------------------------------|-----------|
| 5 | Util | 37 |
| 5.1 | nlpaug.util.file.download | 37 |
| 6 | Indices and tables | 39 |
| | Python Module Index | 41 |
| | Index | 43 |

nlpgaug is a library for textual augmentation in machine learning experiments. The goal is improving deep learning model performance by generating textual data. It also able to generate adversarial examples to prevent adversarial attacks.

CHAPTER 1

Overview

This python library helps you with augmenting nlp for your machine learning projects. Visit [this introduction](#) to understand about Data Augmentation in NLP. Augmenter is the basic element of augmentation while Flow is a pipeline to orchestra multi augmenter together.

- [Data Augmentation library for Text](#)
- [Data Augmentation library for Speech Recognition](#)
- [Data Augmentation library for Audio](#)
- [Does your NLP model able to prevent adversarial attack?](#)

CHAPTER 2

Example

The following examples show a standard use case for augmenter.

- [Audio augmenters](#)
- [Textual augmenters](#)
- [Spectrogram augmenters](#)
- [Custom augmenter](#)
- [TF-IDF model training](#)
- [Flow](#)

3.1 Audio Augmenter

3.1.1 nlpaug.augmenter.audio.crop

Augmenter that apply cropping operation to audio.

```
class nlpaug.augmenter.audio.crop.CropAug(sampling_rate=None, zone=(0.2, 0.8), coverage=0.1, duration=None, crop_range=(0.2, 0.8), crop_factor=2, name='Crop_Aug', verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **sampling_rate** (*int*) – Sampling rate of input audio. Mandatory if duration is provided.
- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 0.1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 25.2 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **duration** (*int*) – Duration of augmentation (in second). Default value is None. If value is provided. *coverage* value will be ignored.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa
>>> aug = naa.CropAug(sampling_rate=44010)
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.2 nlpaug.augmenter.audio.loudness

Augmenter that apply adjusting loudness operation to audio.

```
class nlpaug.augmenter.audio.loudness.LoudnessAug (zone=(0.2, 0.8), coverage=1.0, factor=(0.5, 2), loudness_factor=(0.5, 2), name='Loudness_Aug', verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 42 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **factor** (*tuple*) – Input data volume will be increased (decreased). Augmented value will be picked within the range of this tuple value. Volume will be reduced if value is between 0 and 1.
- **loudness_factor** (*tuple*) – Deprecated. Use *factor* indeed.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa
>>> aug = naa.LoudnessAug()
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation

- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.3 nlpaug.augmenter.audio.mask

Augmenter that apply mask operation to audio.

```
class nlpaug.augmenter.audio.mask.MaskAug (sampling_rate=None, zone=(0.2, 0.8),  
                                           coverage=1.0, duration=(0.2, 0.8),  
                                           mask_range=(0.2, 0.8), mask_factor=2,  
                                           mask_with_noise=True, name='Mask_Aug',  
                                           verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **sampling_rate** (*int*) – Sampling rate of input audio. Mandatory if duration is provided.
- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 42 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **duration** (*float*) – Duration of augmentation (in second). Default value is None. If value is provided. *coverage* value will be ignored.
- **mask_with_noise** (*bool*) – If it is True, targeting area will be replaced by noise. Otherwise, it will be replaced by 0.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa  
>>> aug = naa.MaskAug(sampling_rate=44010)
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.4 nlpaug.augmenter.audio.noise

Augmenter that apply noise injection operation to audio.

```
class nlpaug.augmenter.audio.noise.NoiseAug (zone=(0.2, 0.8), coverage=1.0,
                                             color='white', noises=None,
                                             name='Noise_Aug', verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 42 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **color** (*str*) – Colors of noise. Supported 'white', 'pink', 'red', 'brown', 'brownian', 'blue', 'azure', 'violet', 'purple' and 'random'. If 'random' is used, noise color will be picked randomly in each augment.
- **noises** (*list*) – Background noises for noise injection. You can provide more than one background noise and noise will be picked randomly. Expected format is list of numpy array. If this value is provided. *color* value will be ignored
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa
>>> aug = naa.NoiseAug()
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.5 nlpaug.augmenter.audio.pitch

Augmenter that apply pitch adjustment operation to audio.

```
class nlpaug.augmenter.audio.pitch.PitchAug(sampling_rate, zone=(0.2, 0.8), coverage=1.0, duration=None, factor=(-10, 10), pitch_range=(-10, 10), name='Pitch_Aug', verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **sampling_rate** (*int*) – Sampling rate of input audio. Mandatory if duration is provided.
- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 42 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **duration** (*int*) – Duration of augmentation (in second). Default value is None. If value is provided. *coverage* value will be ignored.
- **pitch_range** (*tuple*) – Deprecated. Use *factor* indeed
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa
>>> aug = naa.PitchAug(sampling_rate=44010)
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.6 nlpaug.augmenter.audio.shift

Augmenter that apply shifting operation to audio.

```
class nlpaug.augmenter.audio.shift.ShiftAug (sampling_rate, duration=3, direction='random', shift_max=3, shift_direction='both', name='Shift_Aug', verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **sampling_rate** (*int*) – Sampling rate of input audio.
- **duration** (*float*) – Max shifting segment (in second)
- **direction** (*str*) – Shifting segment to left, right or one of them. Value can be 'left', 'right' or 'random'
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa
>>> aug = naa.ShiftAug(sampling_rate=44010)
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data


```
>>> augmented_data = aug.augment(data)
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.7 nlpaug.augmenter.audio.speed

Augmenter that apply speed adjustment operation to audio.

```
class nlpaug.augmenter.audio.speed.SpeedAug (zone=(0.2, 0.8), coverage=1.0,
                                             duration=None, factor=(0.5, 2),
                                             speed_range=(0.5, 2), name='Speed_Aug',
                                             verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 42 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **factor** (*int*) – Range of applying speed adjustment operation. Default value is (0.5, 2) Factor for time stretch. Audio will be slowing down if value is between 0 and 1.
- **speed_range** (*tuple*) – Deprecated. Use *factor* indeed
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa
>>> aug = naa.SpeedAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.1.8 nlpaug.augmenter.audio.vtvp

Augmenter that apply vocal tract length perturbation (VTLP) operation to audio.

```
class nlpaug.augmenter.audio.vtvp.VtvpAug (sampling_rate, zone=(0.2, 0.8), coverage=0.1,  
                                           duration=None, fhi=4800, factor=(0.9, 1.1),  
                                           name='Vtvp_Aug', verbose=0)
```

Bases: nlpaug.augmenter.audio.audio_augmenter.AudioAugmenter

Parameters

- **zone** (*tuple*) – Assign a zone for augmentation. Default value is (0.2, 0.8) which means that no any augmentation will be applied in first 20% and last 20% of whole audio.
- **coverage** (*float*) – Portion of augmentation. Value should be between 0 and 1. If 1 is assigned, augment operation will be applied to target audio segment. For example, the audio duration is 60 seconds while zone and coverage are (0.2, 0.8) and 0.7 respectively. 42 seconds $((0.8-0.2)*0.7*60)$ audio will be augmented.
- **factor** (*int*) – Range of applying speed adjustment operation. Default value is (0.5, 2) Factor for time stretch. Audio will be slowing down if value is between 0 and 1.
- **fhi** (*int*) – Boundary frequency. Default value is 4800.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.audio as naa  
>>> aug = naa.VtvpAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.2 Character Augmenter

3.2.1 nlpaug.augmenter.char.keyboard

Augmenter that apply typo error simulation to textual input.

```
class nlpaug.augmenter.char.keyboard.KeyboardAug (name='Keyboard_Aug',
                                                  aug_char_min=1,
                                                  aug_char_max=10,
                                                  aug_char_p=0.3, aug_word_p=0.3,
                                                  aug_word_min=1,
                                                  aug_word_max=10,          stop-
                                                  words=None,          tokenizer=None,
                                                  reverse_tokenizer=None,          in-
                                                  clude_special_char=True,
                                                  include_numeric=True,          in-
                                                  clude_upper_case=True, lang='en',
                                                  verbose=0, stopwords_regex=None,
                                                  model_path=None,          min_char=4,
                                                  include_detail=False)
```

Bases: nlpaug.augmenter.char.char_augmenter.CharAugmenter

Augmenter that simulate typo error by random values. For example, people may type i as o incorrectly. One keyboard distance is leveraged to replace character by possible keyboard error.

Parameters

- **aug_char_p** (*float*) – Percentage of character (per token) will be augmented.
- **aug_char_min** (*int*) – Minimum number of character will be augmented.
- **aug_char_max** (*int*) – Maximum number of character will be augmented. If None is passed, number of augmentation is calculated via `aug_char_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_char_p`. Otherwise, using `aug_max`.
- **aug_word_p** (*float*) – Percentage of word will be augmented.
- **aug_word_min** (*int*) – Minimum number of word will be augmented.
- **aug_word_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via `aug_word_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_word_p`. Otherwise, using `aug_max`.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.

- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_special_char** (*bool*) – Include special character
- **include_upper_case** (*bool*) – If True, upper case character may be included in augmented data.
- **include_numeric** (*bool*) – If True, numeric character may be included in augmented data.
- **min_char** (*int*) – If word less than this value, do not draw word for augmentation
- **model_path** (*str*) – Loading customize model from file system
- **lang** (*str*) – Indicate built-in language model. If custom model is used (passing `model_path`), this value will be ignored.
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.char as nac
>>> aug = nac.KeyboardAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.2.2 nlpaug.augmenter.char.ocr

Augmenter that apply ocr error simulation to textual input.

```
class nlpaug.augmenter.char.ocr.OcrAug (name='OCR_Aug',          aug_char_min=1,
                                       aug_char_max=10,          aug_char_p=0.3,
                                       aug_word_p=0.3,           aug_word_min=1,
                                       aug_word_max=10,           stopwords=None, tok-
                                       enizer=None, reverse_tokenizer=None, ver-
                                       bose=0, stopwords_regex=None, min_char=1,
                                       include_detail=False)
```

Bases: nlpaug.augmenter.char.char_augmenter.CharAugmenter

Augmenter that simulate ocr error by random values. For example, OCR may recognize I as l incorrectly. Pre-defined OCR mapping is leveraged to replace character by possible OCR error.

Parameters

- **aug_char_p** (*float*) – Percentage of character (per token) will be augmented.
- **aug_char_min** (*int*) – Minimum number of character will be augmented.
- **aug_char_max** (*int*) – Maximum number of character will be augmented. If None is passed, number of augmentation is calculated via `aug_char_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_char_p`. Otherwise, using `aug_max`.
- **aug_word_p** (*float*) – Percentage of word will be augmented.
- **aug_word_min** (*int*) – Minimum number of word will be augmented.
- **aug_word_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via `aug_word_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_word_p`. Otherwise, using `aug_max`.
- **min_char** (*int*) – If word less than this value, do not draw word for augmentation
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.char as nac
>>> aug = nac.OcrAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n*=1, *num_thread*=1)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.2.3 nlpaug.augmenter.char.random

Augmenter that apply random character error to textual input.

```
class nlpaug.augmenter.char.random.RandomCharAug (action='substitute',  
                                                name='RandomChar_Aug',  
                                                aug_char_min=1,  
                                                aug_char_max=10,  
                                                aug_char_p=0.3, aug_word_p=0.3,  
                                                aug_word_min=1,  
                                                aug_word_max=10,           in-  
                                                clude_upper_case=True,       in-  
                                                clude_lower_case=True,       in-  
                                                clude_numeric=True, min_char=4,  
                                                swap_mode='adjacent',  
                                                spec_char='!@#$$%^&*()_+',  
                                                stopwords=None, tokenizer=None,  
                                                reverse_tokenizer=None,       ver-  
                                                bose=0,      stopwords_regex=None,  
                                                include_detail=False,      candi-  
                                                ates=None)
```

Bases: nlpaug.augmenter.char.char_augmenter.CharAugmenter

Augmenter that generate character error by random values. For example, people may type i as o incorrectly.

Parameters

- **action** (*str*) – Possible values are 'insert', 'substitute', 'swap' and 'delete'. If value is 'insert', a new character will be injected to randomly. If value is 'substitute', a random character will be replaced original character randomly. If value is 'swap', adjacent characters within sample word will be swapped randomly. If value is 'delete', character will be removed randomly.
- **aug_char_p** (*float*) – Percentage of character (per token) will be augmented.
- **aug_char_min** (*int*) – Minimum number of character will be augmented.
- **aug_char_max** (*int*) – Maximum number of character will be augmented. If None is passed, number of augmentation is calculated via *aug_char_p*. If calculated result from *aug_p* is smaller than *aug_max*, will use calculated result from *aug_char_p*. Otherwise, using *aug_max*.
- **aug_word_p** (*float*) – Percentage of word will be augmented.
- **aug_word_min** (*int*) – Minimum number of word will be augmented.

- **aug_word_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via `aug_word_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_word_p`. Otherwise, using `aug_max`.
- **include_upper_case** (*bool*) – If True, upper case character may be included in augmented data. If ‘candidates’ value is provided, this param will be ignored.
- **include_lower_case** (*bool*) – If True, lower case character may be included in augmented data. If ‘candidates’ value is provided, this param will be ignored.
- **include_numeric** (*bool*) – If True, numeric character may be included in augmented data. If ‘candidates’ value is provided, this param will be ignored.
- **min_char** (*int*) – If word less than this value, do not draw word for augmentation
- **swap_mode** – When action is ‘swap’, you may pass ‘adjacent’, ‘middle’ or ‘random’. ‘adjacent’ means swap action only consider adjacent character (within same word). ‘middle’ means swap action consider adjacent character but not the first and last character of word. ‘random’ means swap action will be executed without constraint.
- **spec_char** (*str*) – Special character may be included in augmented data. If ‘candidates’ value is provided, this param will be ignored.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **candidates** (*List*) – List of string for augmentation. E.g. [‘AAA’, ‘11’, ‘===’]. If values is provided, `include_upper_case`, `include_lower_case`, `include_numeric` and `spec_char` will be ignored.
- **name** (*str*) – Name of this augmenter.

```
>>> import nlpaug.augmenter.char as nac
>>> aug = nac.RandomCharAug()
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data

- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.3 Sentence Augmenter

3.3.1 nlpaug.augmenter.sentence.context_word_embs_sentence

Augmenter that apply operation (sentence level) to textual input based on contextual word embeddings.

class nlpaug.augmenter.sentence.context_word_embs_sentence.ContextualWordEmbsForSentenceAugmenter

Bases: nlpaug.augmenter.sentence.sentence_augmenter.SentenceAugmenter

Augmenter that leverage contextual word embeddings to find top n similar word for augmentation.

Parameters

- **model_path** (*str*) – Model name or model path. It used transformers to load the model. Tested 'xlnet-base-cased', 'gpt2', 'distilgpt2'. If you want to reduce inference time, you may select *distilgpt2*.
- **temperature** (*float*) – Controlling randomness. Default value is 1 and lower temperature results in less random behavior
- **top_k** (*int*) – Controlling lucky draw pool. Top k score token will be used for augmentation. Larger k, more token can be used. Default value is 100. If value is None which means using all possible tokens.
- **top_p** (*float*) – Controlling lucky draw pool. Top p of cumulative probability will be removed. Larger p, more token can be used. Default value is None which means using all possible tokens.
- **device** (*str*) – Use either cpu or gpu. Default value is None, it uses GPU if having. While possible values are 'cuda' and 'cpu'.

- **force_reload** (*bool*) – Force reload the contextual word embeddings model to memory when initialize the class. Default value is False and suggesting to keep it as False if performance is the consideration.
- **optimize** (*obj*) – Configuration for optimized process. *external_memory*: Persisting previous computed result for next prediction. Extra memory will be used in order to have shorter inference time. *gpt2* and ‘*distilgpt2*’ are supported.
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.sentence as nas
>>> aug = nas.ContextualWordEmbsForSentenceAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.4 Spectrogram Augmenter

3.4.1 nlpaug.augmenter.spectrogram.frequency_masking

Augmenter that apply frequency based masking to spectrogram input.

```
class nlpaug.augmenter.spectrogram.frequency_masking.FrequencyMaskingAug (mask_factor,
                                                                    name='FrequencyMaskingAug',
                                                                    verbose=0)
```

Bases: `nlpaug.augmenter.spectrogram.spectrogram_augmenter.SpectrogramAugmenter`

Augmenter that mask spectrogram based on frequency by random values.

Parameters

- **mask_factor** (*int*) – Value between 0 and mask_factor will be picked randomly. Mask range will be between [0, v - master_factor) while v is the number of mel frequency channels.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.spectrogram as nas
>>> aug = nas.FrequencyMaskingAug(mask_factor=80)
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.4.2 nlpaug.augmenter.spectrogram.time_masking

Augmenter that apply time based masking to spectrogram input.

```
class nlpaug.augmenter.spectrogram.time_masking.TimeMaskingAug(mask_factor,
                                                                name='TimeMasking_Aug',
                                                                verbose=0)
```

Bases: nlpaug.augmenter.spectrogram.spectrogram_augmenter.
SpectrogramAugmenter

Augmenter that mask spectrogram based on time by random values.

Parameters

- **mask_factor** (*int*) – Value between 0 and mask_factor will be picked randomly. Mask range will be between [0, tau - master_factor) while tau is time range of input.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.spectrogram as nas
>>> aug = nas.TimeMaskingAug(mask_factor=80)
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augment_s (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5 Word Augmenter

3.5.1 nlpaug.augmenter.word.antonym

Augmenter that apply semantic meaning based to textual input.

```
class nlpaug.augmenter.word.antonym.AntonymAug (name='Antonym_Aug',      aug_min=1,
                                                  aug_max=10,      aug_p=0.3,
                                                  lang='eng',    stopwords=None, tok-
                                                  enizer=None, reverse_tokenizer=None,
                                                  stopwords_regex=None,      in-
                                                  clude_detail=False, verbose=0)
```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that leverage semantic meaning to substitute word.

Parameters

- **lang** (*str*) – Language of your text. Default value is 'eng'.
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via aug_p. If calculated result from aug_p is smaller than aug_max, will use calculated result from aug_p. Otherwise, using aug_max.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.

- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.AntonymAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5.2 nlpaug.augmenter.word.context_word_embs

Augmenter that apply operation (word level) to textual input based on contextual word embeddings.

```

class nlpaug.augmenter.word.context_word_embs.ContextualWordEmbsAug(model_path='bert-
base-uncased',
action='substitute',
temperature=1.0,
top_k=100,
top_p=None,
name='ContextualWordEmbs_Au
aug_min=1,
aug_max=10,
aug_p=0.3,
stop-
words=None,
skip_unknown_word=False,
de-
vice=None,
force_reload=False,
opti-
mize=None,
stop-
words_regex=None,
ver-
bose=0,
in-
clude_detail=False)

```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that leverage contextual word embeddings to find top n similar word for augmentation.

Parameters

- **model_path** (*str*) – Model name or model path. It used transformers to load the model. Tested ‘bert-base-uncased’, ‘bert-base-cased’, ‘distilbert-base-uncased’, ‘roberta-base’, ‘distilroberta-base’, ‘xlnet-base-cased’.
- **action** (*str*) – Either ‘insert’ or ‘substitute’. If value is ‘insert’, a new word will be injected to random position according to contextual word embeddings calculation. If value is ‘substitute’, word will be replaced according to contextual embeddings calculation
- **temperature** (*float*) – Controlling randomness. Default value is 1 and lower temperature results in less random behavior
- **top_k** (*int*) – Controlling lucky draw pool. Top k score token will be used for augmentation. Larger k, more token can be used. Default value is 100. If value is None which means using all possible tokens.
- **top_p** (*float*) – Controlling lucky draw pool. Top p of cumulative probability will be removed. Larger p, more token can be used. Default value is None which means using all possible tokens.
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via aug_p. If calculated result from aug_p is smaller than aug_max, will use calculated result from aug_p. Otherwise, using aug_max.

- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **skip_unknown_word** (*bool*) – Do not substitute unknown word (e.g. AAAAAAAAAA)
- **device** (*str*) – Use either cpu or gpu. Default value is None, it uses GPU if having. While possible values are 'cuda' and 'cpu'.
- **force_reload** (*bool*) – Force reload the contextual word embeddings model to memory when initialize the class. Default value is False and suggesting to keep it as False if performance is the consideration.
- **optimize** (*bool*) – If true, optimized process will be executed. For example, GPT2 will use "return_past" to reduce inference time.
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augementer

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.ContextualWordEmbsAug()
```

augment (*data, n=1, num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data, n=1, num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

device = None

TODO: Reserve 2 spaces (e.g. [CLS], [SEP]) is not enough as it hit CUDA error in batch processing mode. Therefore, forcing to reserve 5 times of reserved spaces (i.e. 5)

3.5.3 nlpaug.augmenter.word.random

Augmenter that apply random word operation to textual input.

```
class nlpaug.augmenter.word.random.RandomWordAug (action='delete',
                                                  name='RandomWord_Aug',
                                                  aug_min=1,          aug_max=10,
                                                  aug_p=0.3, stopwords=None, tar-
                                                  get_words=None, tokenizer=None,
                                                  reverse_tokenizer=None,
                                                  stopwords_regex=None,      in-
                                                  clude_detail=False, verbose=0)
```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that apply randomly behavior for augmentation.

Parameters

- **action** (*str*) – ‘substitute’, ‘swap’, ‘delete’ or ‘crop’. If value is ‘swap’, adjacent words will be swapped randomly. If value is ‘delete’, word will be removed randomly. If value is ‘crop’, a set of continuous word will be removed randomly.
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via aug_p. If calculated result from aug_p is smaller than aug_max, will use calculated result from aug_p. Otherwise, using aug_max.
- **stopwords** (*list*) – List of words which will be skipped from augment operation. Not effective if action is ‘crop’
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation. Not effective if action is ‘crop’
- **target_words** (*list*) – List of word for replacement (used for substitute operation only). Default value is _.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.RandomWordAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5.4 nlpaug.augmenter.word.spelling

Augmenter that apply spelling error simulation to textual input.

```
class nlpaug.augmenter.word.spelling.SpellingAug(dict_path,      name='Spelling_Aug',
                                                  aug_min=1,        aug_max=10,
                                                  aug_p=0.3,        stopwords=None,
                                                  tokenizer=None,      re-
                                                  verse_tokenizer=None, in-
                                                  clude_reverse=True,  stop-
                                                  words_regex=None,   in-
                                                  clude_detail=False, verbose=0)
```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that leverage pre-defined spelling mistake dictionary to simulate spelling mistake.

Parameters

- **dict_path** (*str*) – Path of misspelling dictionary
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via `aug_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_p`. Otherwise, using `aug_max`.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.SpellingAug(dict_path='./spelling_en.txt')
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output

- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5.5 nlpaug.augmenter.word.split

Augmenter that apply word splitting operation to textual input.

```
class nlpaug.augmenter.word.split.SplitAug (name='Split_Aug',          aug_min=1,
                                           aug_max=10, aug_p=0.3, min_char=4,
                                           stopwords=None, tokenizer=None,
                                           reverse_tokenizer=None, stop-
                                           words_regex=None, include_detail=False,
                                           verbose=0)
```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that apply word splitting for augmentation.

Parameters

- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via `aug_p`. If calculated result from `aug_p` is smaller than `aug_max`, will use calculated result from `aug_p`. Otherwise, using `aug_max`.
- **min_char** (*int*) – If word less than this value, do not draw word for augmentation
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.SplitAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5.6 nlpaug.augmenter.word.synonym

Augmenter that apply semantic meaning based to textual input.

```
class nlpaug.augmenter.word.synonym.SynonymAug (aug_src='wordnet', model_path=None,
                                                name='Synonym_Aug',    aug_min=1,
                                                aug_max=10,             aug_p=0.3,
                                                lang='eng', stopwords=None, tokenizer=None, reverse_tokenizer=None,
                                                stopwords_regex=None,
                                                force_reload=False,      in-
                                                clude_detail=False, verbose=0)
```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that leverage semantic meaning to substitute word.

Parameters

- **aug_src** (*str*) – Support ‘wordnet’ and ‘ppdb’ .
- **model_path** (*str*) – Path of dictionary. Mandatory field if using PPDB as data source
- **lang** (*str*) – Language of your text. Default value is ‘eng’.
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via aug_p. If calculated result from aug_p is smaller than aug_max, will use calculated result from aug_p. Otherwise, using aug_max.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.

- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **force_reload** (*bool*) – Force reload model to memory when initialize the class. Default value is False and suggesting to keep it as False if performance is the consideration.
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.SynonymAug()
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5.7 nlpaug.augmenter.word.tfidf

Augmenter that apply TF-IDF based to textual input.

```
class nlpaug.augmenter.word.tfidf.TfIdfAug(model_path='', action='substitute',
                                           name='Tfidf_Aug', aug_min=1,
                                           aug_max=10, aug_p=0.3, top_k=5,
                                           stopwords=None, tokenizer=None,
                                           reverse_tokenizer=None, stop-
                                           words_regex=None, include_detail=False,
                                           verbose=0)
```

Bases: `nlpaug.augmenter.word.word_augmenter.WordAugmenter`

Augmenter that leverage TF-IDF statistics to insert or substitute word.

Parameters

- **model_path** (*str*) – Downloaded model directory. Either model_path or model is must be provided
- **action** (*str*) – Either ‘insert or ‘substitute’. If value is ‘insert’, a new word will be injected to random position according to TF-IDF calculation. If value is ‘substitute’, word will be replaced according to TF-IDF calculation
- **top_k** (*int*) – Controlling lucky draw pool. Top k score token will be used for augmentation. Larger k, more token can be used. Default value is 5. If value is None which means using all possible tokens.
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via aug_p. If calculated result from aug_p is smaller than aug_max, will use calculated result from aug_p. Otherwise, using aug_max.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.TfIdfAug(model_path='.')
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

3.5.8 nlpaug.augmenter.word.word_embs

Augmenter that apply operation to textual input based on word embeddings.

```
class nlpaug.augmenter.word.word_embs.WordEmbsAug(model_type, model_path='',
                                                    model=None, action='substitute',
                                                    name='WordEmbs_Aug',
                                                    aug_min=1, aug_max=10,
                                                    aug_p=0.3, top_k=100,
                                                    n_gram_separator='_', stopwords=None,
                                                    tokenizer=None,
                                                    reverse_tokenizer=None,
                                                    force_reload=False, stopwords_regex=None,
                                                    include_detail=False, verbose=0)
```

Bases: nlpaug.augmenter.word.word_augmenter.WordAugmenter

Augmenter that leverage word embeddings to find top n similar word for augmentation.

Parameters

- **model_type** (*str*) – Model type of word embeddings. Expected values include 'word2vec', 'glove' and 'fasttext'.
- **model_path** (*str*) – Downloaded model directory. Either model_path or model is must be provided
- **model** (*obj*) – Pre-loaded model
- **action** (*str*) – Either 'insert' or 'substitute'. If value is 'insert', a new word will be injected to random position according to word embeddings calculation. If value is 'substitute', word will be replaced according to word embeddings calculation
- **top_k** (*int*) – Controlling lucky draw pool. Top k score token will be used for augmentation. Larger k, more token can be used. Default value is 100. If value is None which means using all possible tokens.
- **aug_p** (*float*) – Percentage of word will be augmented.
- **aug_min** (*int*) – Minimum number of word will be augmented.
- **aug_max** (*int*) – Maximum number of word will be augmented. If None is passed, number of augmentation is calculated via aug_p. If calculated result from aug_p is smaller than aug_max, will use calculated result from aug_p. Otherwise, using aug_max.
- **stopwords** (*list*) – List of words which will be skipped from augment operation.
- **stopwords_regex** (*str*) – Regular expression for matching words which will be skipped from augment operation.
- **tokenizer** (*func*) – Customize tokenization process
- **reverse_tokenizer** (*func*) – Customize reverse of tokenization process
- **force_reload** (*bool*) – If True, model will be loaded every time while it takes longer time for initialization.
- **include_detail** (*bool*) – Change detail will be returned if it is True.
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.augmenter.word as naw
>>> aug = naw.WordEmbsAug(model_type='word2vec', model_path='.')
```

augment (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*object*) – Data for augmentation
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1

Returns Augmented data

```
>>> augmented_data = aug.augment(data)
```

augments (*data*, *n=1*, *num_thread=1*)

Parameters

- **data** (*list*) – List of data
- **n** (*int*) – Number of unique augmented output
- **num_thread** (*int*) – Number of thread for data augmentation. Use this option when you are using CPU and n is larger than 1. Do NOT support GPU process.

Returns Augmented data (Does not follow original order)

```
>>> augmented_data = aug.augment(data)
```

4.1 nlpaug.flow.sequential

Flow that apply augmentation sequentially.

```
class nlpaug.flow.sequential.Sequential (flow=None, name='Sequential_Pipeline', include_detail=False, verbose=0)
```

Bases: nlpaug.flow.pipeline.Pipeline

Flow that apply augmenters sequentially.

Parameters

- **flow** (*list*) – list of flow or augmenter
- **name** (*str*) – Name of this augmenter

```
>>> import nlpaug.flow as naf
>>> import nlpaug.augmenter.char as nac
>>> import nlpaug.augmenter.word as naw
>>> flow = naf.Sequential([nac.RandomCharAug(), naw.RandomWordAug()])
```

4.2 nlpaug.flow.sometimes

Flow that apply augmentation randomly.

```
class nlpaug.flow.sometimes.Sometimes (flow=None, name='Sometimes_Pipeline', pipeline_p=0.2, aug_p=1, include_detail=False, verbose=0)
```

Bases: nlpaug.flow.pipeline.Pipeline

Flow that apply augmenters randomly.

Parameters

- **flow** (*list*) – list of flow or augmenter

- **name** (*str*) – Name of this augementer

```
>>> import nlpaug.flow as naf
>>> import nlpaug.augmenter.char as nac
>>> import nlpaug.augmenter.word as naw
>>> flow = naf.Sometimes([nac.RandomCharAug(), naw.RandomWordAug()])
```


5.1 nlpaug.util.file.download

class nlpaug.util.file.download.DownloadUtil

Bases: object

Helper function for downloading external dependency

```
>>> from nlpaug.util.file.download import DownloadUtil
```

static download_fasttext(*model_name*, *dest_dir*)

Parameters

- **model_name** (*str*) – GloVe pre-trained model name. Possible values are ‘wiki-news-300d-1M’, ‘wiki-news-300d-1M-subword’, ‘crawl-300d-2M’ and ‘crawl-300d-2M-subword’
- **dest_dir** (*str*) – Directory of saving file

```
>>> DownloadUtil.download_fasttext('glove.6B', '.')
```

static download_glove(*model_name*, *dest_dir*)

Parameters

- **model_name** (*str*) – GloVe pre-trained model name. Possible values are ‘glove.6B’, ‘glove.42B.300d’, ‘glove.840B.300d’ and ‘glove.twitter.27B’
- **dest_dir** (*str*) – Directory of saving file

```
>>> DownloadUtil.download_glove('glove.6B', '.')
```

static download_word2vec(*dest_dir*)

Parameters **dest_dir** (*str*) – Directory of saving file

```
>>> DownloadUtil.download_word2vec('.')
```

static unzip (*file_path*, *dest_dir=None*)

Parameters **file_path** (*str*) – File path for unzip

```
>>> DownloadUtil.unzip('zip_file.zip')
```

See modindex for API.

CHAPTER 6

Indices and tables

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

n

- `nlpaug.augmenter.audio.crop`, 7
- `nlpaug.augmenter.audio.loudness`, 8
- `nlpaug.augmenter.audio.mask`, 9
- `nlpaug.augmenter.audio.noise`, 10
- `nlpaug.augmenter.audio.pitch`, 11
- `nlpaug.augmenter.audio.shift`, 12
- `nlpaug.augmenter.audio.speed`, 13
- `nlpaug.augmenter.audio.vt1p`, 14
- `nlpaug.augmenter.char.keyboard`, 15
- `nlpaug.augmenter.char.ocr`, 16
- `nlpaug.augmenter.char.random`, 18
- `nlpaug.augmenter.sentence.context_word_embs_sentence`,
20
- `nlpaug.augmenter.spectrogram.frequency_masking`,
21
- `nlpaug.augmenter.spectrogram.time_masking`,
22
- `nlpaug.augmenter.word.antonym`, 23
- `nlpaug.augmenter.word.context_word_embs`,
24
- `nlpaug.augmenter.word.random`, 26
- `nlpaug.augmenter.word.spelling`, 28
- `nlpaug.augmenter.word.split`, 29
- `nlpaug.augmenter.word.synonym`, 30
- `nlpaug.augmenter.word.tfidf`, 31
- `nlpaug.augmenter.word.word_embs`, 33
- `nlpaug.flow.sequential`, 35
- `nlpaug.flow.sometimes`, 35
- `nlpaug.util.file.download`, 37

A

- AntonymAug (class in `nlpaug.augmenter.word.antonym`), 23
- `augment()` (`nlpaug.augmenter.audio.crop.CropAug` method), 7
- `augment()` (`nlpaug.augmenter.audio.loudness.LoudnessAug` method), 8
- `augment()` (`nlpaug.augmenter.audio.mask.MaskAug` method), 9
- `augment()` (`nlpaug.augmenter.audio.noise.NoiseAug` method), 10
- `augment()` (`nlpaug.augmenter.audio.pitch.PitchAug` method), 11
- `augment()` (`nlpaug.augmenter.audio.shift.ShiftAug` method), 12
- `augment()` (`nlpaug.augmenter.audio.speed.SpeedAug` method), 13
- `augment()` (`nlpaug.augmenter.audio.vtlf.VtlfAug` method), 14
- `augment()` (`nlpaug.augmenter.char.keyboard.KeyboardAug` method), 16
- `augment()` (`nlpaug.augmenter.char.ocr.OcrAug` method), 17
- `augment()` (`nlpaug.augmenter.char.random.RandomCharAug` method), 19
- `augment()` (`nlpaug.augmenter.sentence.context_word_embs_sentence.ContextualWordEmbsForSentenceAug` method), 21
- `augment()` (`nlpaug.augmenter.spectrogram.frequency_masking.FrequencyMaskingAug` method), 22
- `augment()` (`nlpaug.augmenter.spectrogram.time_masking.TimeMaskingAug` method), 22
- `augment()` (`nlpaug.augmenter.word.antonym.AntonymAug` method), 24
- `augment()` (`nlpaug.augmenter.word.context_word_embs.ContextualWordEmbsAug` method), 26
- `augment()` (`nlpaug.augmenter.word.random.RandomWordAug` method), 27
- `augment()` (`nlpaug.augmenter.word.spelling.SpellingAug` method), 28
- `augment()` (`nlpaug.augmenter.word.split.SplitAug` method), 29
- `augment()` (`nlpaug.augmenter.word.synonym.SynonymAug` method), 31
- `augment()` (`nlpaug.augmenter.word.tfidf.TfidfAug` method), 32
- `augment()` (`nlpaug.augmenter.word.word_embs.WordEmbsAug` method), 33
- `augments()` (`nlpaug.augmenter.audio.crop.CropAug` method), 8
- `augments()` (`nlpaug.augmenter.audio.loudness.LoudnessAug` method), 9
- `augments()` (`nlpaug.augmenter.audio.mask.MaskAug` method), 10
- `augments()` (`nlpaug.augmenter.audio.noise.NoiseAug` method), 11
- `augments()` (`nlpaug.augmenter.audio.pitch.PitchAug` method), 12
- `augments()` (`nlpaug.augmenter.audio.shift.ShiftAug` method), 13
- `augments()` (`nlpaug.augmenter.audio.speed.SpeedAug` method), 13
- `augments()` (`nlpaug.augmenter.audio.vtlf.VtlfAug` method), 14
- `augments()` (`nlpaug.augmenter.char.keyboard.KeyboardAug` method), 16
- `augments()` (`nlpaug.augmenter.char.ocr.OcrAug` method), 17
- `augments()` (`nlpaug.augmenter.char.random.RandomCharAug` method), 19
- `augments()` (`nlpaug.augmenter.sentence.context_word_embs_sentence.ContextualWordEmbsForSentenceAug` method), 21
- `augments()` (`nlpaug.augmenter.spectrogram.frequency_masking.FrequencyMaskingAug` method), 22
- `augments()` (`nlpaug.augmenter.spectrogram.time_masking.TimeMaskingAug` method), 23
- `augments()` (`nlpaug.augmenter.word.antonym.AntonymAug` method), 24
- `augments()` (`nlpaug.augmenter.word.context_word_embs.ContextualWordEmbsAug` method), 26

augments () (*nlpaug.augmenter.word.random.RandomWordAug*
method), 27

augments () (*nlpaug.augmenter.word.spelling.SpellingAug*
method), 29

augments () (*nlpaug.augmenter.word.split.SplitAug*
method), 30

augments () (*nlpaug.augmenter.word.synonym.SynonymAug*
method), 31

augments () (*nlpaug.augmenter.word.tfidf.TfidfAug*
method), 32

augments () (*nlpaug.augmenter.word.word_embs.WordEmbsAug*
method), 34

C

ContextualWordEmbsAug (class in *nl-
paug.augmenter.word.context_word_embs*),
24

ContextualWordEmbsForSentenceAug
(class in *nl-
paug.augmenter.sentence.context_word_embs_sentence*),
20

CropAug (class in *nlpaug.augmenter.audio.crop*), 7

D

device (*nlpaug.augmenter.word.context_word_embs.ContextualWordEmbsAug*
attribute), 26

download_fasttext () (*nl-
paug.util.file.download.DownloadUtil*
static
method), 37

download_glove () (*nl-
paug.util.file.download.DownloadUtil*
static
method), 37

download_word2vec () (*nl-
paug.util.file.download.DownloadUtil*
static
method), 37

DownloadUtil (class in *nlpaug.util.file.download*), 37

F

FrequencyMaskingAug (class in *nl-
paug.augmenter.spectrogram.frequency_masking*),
21

K

KeyboardAug (class in *nl-
paug.augmenter.char.keyboard*), 15

L

LoudnessAug (class in *nl-
paug.augmenter.audio.loudness*), 8

M

MaskAug (class in *nlpaug.augmenter.audio.mask*), 9

N

nlpaug.augmenter.audio.crop (module), 7

nlpaug.augmenter.audio.loudness (module),
8

nlpaug.augmenter.audio.mask (module), 9

nlpaug.augmenter.audio.noise (module), 10

nlpaug.augmenter.audio.pitch (module), 11

nlpaug.augmenter.audio.shift (module), 12

nlpaug.augmenter.audio.speed (module), 13

nlpaug.augmenter.audio.vt1p (module), 14

nlpaug.augmenter.char.keyboard (module),
15

nlpaug.augmenter.char.ocr (module), 16

nlpaug.augmenter.char.random (module), 18

nlpaug.augmenter.sentence.context_word_embs_sentence
(module), 20

nlpaug.augmenter.spectrogram.frequency_masking
(module), 21

nlpaug.augmenter.spectrogram.time_masking
(module), 22

nlpaug.augmenter.word.antonym (module), 23

nlpaug.augmenter.word.context_word_embs
(module), 24

nlpaug.augmenter.word.random (module), 26

nlpaug.augmenter.word.spelling (module),
28

nlpaug.augmenter.word.split (module), 29

nlpaug.augmenter.word.synonym (module), 30

nlpaug.augmenter.word.tfidf (module), 31

nlpaug.augmenter.word.word_embs (module),
33

nlpaug.flow.sequential (module), 35

nlpaug.flow.sometimes (module), 35

nlpaug.util.file.download (module), 37

NoiseAug (class in *nlpaug.augmenter.audio.noise*), 10

O

OcrAug (class in *nlpaug.augmenter.char.ocr*), 16

P

PitchAug (class in *nlpaug.augmenter.audio.pitch*), 11

R

RandomCharAug (class in *nl-
paug.augmenter.char.random*), 18

RandomWordAug (class in *nl-
paug.augmenter.word.random*), 26

S

Sequential (class in *nlpaug.flow.sequential*), 35

ShiftAug (class in *nlpaug.augmenter.audio.shift*), 12

Sometimes (class in *nlpaug.flow.sometimes*), 35

SpeedAug (class in *nlpaug.augmenter.audio.speed*), 13

SpellingAug (class in nl-
paug.augmenter.word.spelling), [28](#)

SplitAug (class in nlpaug.augmenter.word.split), [29](#)

SynonymAug (class in nl-
paug.augmenter.word.synonym), [30](#)

T

TfIdfAug (class in nlpaug.augmenter.word.tfidf), [31](#)

TimeMaskingAug (class in nl-
paug.augmenter.spectrogram.time_masking),
[22](#)

U

unzip() (nlpaug.util.file.download.DownloadUtil
static method), [38](#)

V

Vt1pAug (class in nlpaug.augmenter.audio.vt1p), [14](#)

W

WordEmbsAug (class in nl-
paug.augmenter.word.word_embs), [33](#)