
Under The Sea Documentation

Release 1.1.4

Vu Anh

Dec 30, 2017

Notes

1 Underthesea - Vietnamese NLP Toolkit	3
2 Credits	7
3 History	9
4 word_sent	11
5 pos_tag	13
6 chunking	15
7 ner	17
8 classify	19
9 sentiment	21
10 Indices and tables	23
Python Module Index	25

Vietnamese NLP Toolkit

CHAPTER 1

Underthesea - Vietnamese NLP Toolkit

Documentation: <https://underthesea.readthedocs.io>

Live demo: [underthesea app](#)

Facebook Page: <https://www.facebook.com/undertheseanlp/>

1.1 Installation

To install underthesea, simply:

```
$ pip install underthesea==1.1.6rc2
```

Satisfaction, guaranteed.

1.2 Usage

- 1. *Word Segmentation*
- 2. *POS Tagging*
- 3. *Chunking*
- 4. *Named Entity Recognition*
- 5. *Text Classification*
- 6. *Sentiment Analysis*

1.2.1 1. Word Segmentation

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import word_sent
>>> sentence = u"Chúng ta thng nói dn Rau sch, Rau an toàn đ phân bit vi các rau bình_
→thng bán ngoài ch."
>>> word_sent(sentence)
[u"Chúng ta", u"thng", u"nói", u"dn", u"Rau sch", u",", u"Rau", u"an toàn", u"đ", u
→"phân bit", u"vi",
u"các", u"rau", u"bình thng", u"bán", u"ngoài", u"ch", u"."]
>>> word_sent(sentence, format="text")
u'Chúng_ta thng nói dn Rau_sch , Rau_an_toàn đ phân_bit vi các rau_bình_thng bán_
→ngoài ch .'
```

1.2.2 2. POS Tagging

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import pos_tag
>>> text = u"Ch tht chó ni ting TP H Chí Minh b truy quét"
>>> pos_tag(text)
[(u'Ch', 'N'),
(u'tht', 'N'),
(u'chó', 'N'),
(u'ni ting', 'A'),
(u'', 'E'),
(u'TP HCM', 'Np'),
(u'b', 'V'),
(u'truy quét', 'V')]
```

1.2.3 3. Chunking

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import chunk
>>> text = u"Bác sĩ bây gi có th thn nhién báo tin bnh nhân b ung th?"
>>> chunk(text)
[(u'Bác sĩ', 'N', 'B-NP'),
(u'bây gi', 'P', 'I-NP'),
(u'có th', 'R', 'B-VP'),
(u'thn nhién', 'V', 'I-VP'),
(u'báo tin', 'N', 'B-NP'),
(u'bnh nhân', 'N', 'I-NP'),
(u'b', 'V', 'B-VP'),
(u'ung th', 'N', 'I-VP'),
(u'?', 'CH', 'O')]
```

1.2.4 4. Named Entity Recognition

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import ner
>>> text = u"Cha tit l lch trình ti Vit Nam ca Tng thng M Donald Trump"
```

```
>>> ner(text)
[('Cha', 'R', 'O', 'O'),
 ('tit l', 'V', 'B-VP', 'O'),
 ('lch trình', 'V', 'B-VP', 'O'),
 ('ti', 'E', 'B-PP', 'O'),
 ('Vit Nam', 'Np', 'B-NP', 'B-LOC'),
 ('ca', 'E', 'B-PP', 'O'),
 ('Tng thng', 'N', 'B-NP', 'O'),
 ('M', 'Np', 'B-NP', 'B-LOC'),
 ('Donald', 'Np', 'B-NP', 'B-PER'),
 ('Trump', 'Np', 'B-NP', 'I-PER')]
```

1.2.5 5. Text Classification

```
$ pip install Cython
$ pip install future scipy numpy scikit-learn
$ pip install -U fasttext --no-cache-dir --no-deps --force-reinstall
$ underthesea data
```

Usage

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import classify
>>> classify("HLV đư tiên Premier League b sa thi sau 4 vòng đư")
['The thao']
>>> classify("Hi đng t vn kinh doanh Asean vinh danh gii thng quc t")
['Kinh doanh']
>>> classify("Đánh giá "rp hát ti gia" Samsung Soundbar Sound+ MS750")
['Vi tinh']
```

1.2.6 6. Sentiment Analysis

```
$ pip install future scipy numpy scikit-learn==0.19.0 joblib
```

Usage

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import sentiment
>>> sentiment("Gi my ln mà lúc nào cũng là các chuyên viên đang bn ht ")
('CUSTOMER SUPPORT#NEGATIVE',)
>>> sentiment("bidv cho vay hay ko phu thuoc y thich cua thang tham dinh, ko co quy_
↳dinh ro rang")
('LOAN#NEGATIVE',)
```

1.3 Up Coming Features

- Text to Speech
- Automatic Speech Recognition
- Machine Translation
- Dependency Parsing

1.4 Contributing

Do you want to contribute with underthesea development? Great! Please read more details at [CONTRIBUTING.rst](#).

CHAPTER 2

Credits

2.1 Development Lead

- Vu Anh <brother.rain.1024@gmail.com>

2.2 Contributors

- Bui Nhat Anh <buiinhatanh1208@gmail.com>
- Doan Viet Dung <doanvietdung273@gmail.com>

CHAPTER 3

History

3.1 1.1.6 (2017-12-26)

- New feature: aspect sentiment analysis
- Integrate with languageflow 1.1.6
- Fix bug tokenize string with '=' (#159)

3.2 1.1.5 (2017-10-12)

- New feature: named entity recognition
- Refactor and update model for word_sent, pos_tag, chunking

3.3 1.1.4 (2017-09-12)

- New feature: text classification
- [bug] Fix Text error
- [doc] Add facebook link

3.4 1.1.3 (2017-08-30)

- Add live demo: <https://underthesea.herokuapp.com/>

3.5 1.1.2 (2017-08-22)

- Add dictionary

3.6 1.1.1 (2017-07-05)

- Support Python 3
- Refactor feature_engineering code

3.7 1.1.0 (2017-05-30)

- Add chunking feature
- Add pos_tag feature
- Add word_sent feature, fix performance
- Add Corpus class
- Add Transformer classes
- Integrated with dictionary of Ho Ngoc Duc
- Add travis-CI, auto build with PyPI

3.8 1.0.0 (2017-03-01)

- First release on PyPI.
- First release on Readthedocs

CHAPTER 4

word_sent

`underthesea.word_sent.word_sent(sentence, format=None)`

Vietnamese word segmentation

Parameters `sentence` (`{unicode, str}`) – raw sentence

Returns `tokens` – tagged sentence

Return type list of text

Examples

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import word_sent
>>> sentence = "Bác sĩ bây giờ có thể tham gia sự kiện báo tin ban nhân bưng trà"
```

```
>>> word_sent(sentence)
['Bác', 'sĩ', 'bây', 'gi', 'có', 'thể', 'tham', 'gia', 'sự', 'kiện', 'báo', 'tin', 'ban', 'nhân', 'bưng', 'trà']
```

```
>>> word_sent(sentence, format="text")
'Bác_sĩ_bây_gi_có_th_tham_gia_sự_kiện_báo_tin_ban_nhân_b_ung_trà'
```


CHAPTER 5

pos_tag

`underthesea.pos_tag.pos_tag(sentence, format=None)`

Vietnamese POS tagging

Parameters `sentence` (`{unicode, str}`) – Raw sentence

Returns `tokens` – tagged sentence

Return type list of tuple with word, pos tag

Examples

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import pos_tag
>>> sentence = "Ch ht chó ni ting TPHCM b truy quét"
>>> pos_tag(sentence)
[('Ch', 'N'),
 ('ht', 'N'),
 ('chó', 'N'),
 ('ni ting', 'A'),
 ('', 'E'),
 ('TPHCM', 'Np'),
 ('b', 'V'),
 ('truy quét', 'V')]
```


CHAPTER 6

chunking

`underthesea.chunking.chunk(sentence, format=None)`

Vietnamese chunking

Parameters `sentence` (`{unicode, str}`) – raw sentence

Returns `tokens` – tagged sentence

Return type list of tuple with word, pos tag, chunking tag

Examples

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import chunk
>>> sentence = "Nghi vn 4 thi th Triu Tiêu trôi dt b bin Nht Bn"
>>> chunk(sentence)
[('Nghi', 'N', 'B-NP'),
 ('vn', 'N', 'B-NP'),
 ('4', 'M', 'B-NP'),
 ('thi', 'N', 'B-NP'),
 ('th', 'N', 'B-NP'),
 ('Triu', 'Np', 'B-NP'),
 ('Tiêu', 'Np', 'B-NP'),
 ('trôi', 'V', 'B-VP'),
 ('dt', 'DT', 'B-VP'),
 ('b', 'N', 'B-NP'),
 ('bin', 'N', 'B-NP'),
 ('Nht', 'N', 'B-NP'),
 ('Bn', 'N', 'B-NP')]
```


CHAPTER 7

ner

`underthesea.ner.ner(sentence,format=None)`

Location and classify named entities in text

Parameters `sentence` (`{unicode, str}`) – raw sentence

Returns `tokens` – tagged sentence

Return type list of tuple with word, pos tag, chunking tag, ner tag

Examples

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import ner
>>> sentence = "Ông Putin ca ngi nhng thành tu vĩ di ca Liên Xô"
>>> ner(sentence)
[('Ông', 'Nc', 'B-NP', 'O'),
 ('Putin', 'Np', 'B-NP', 'B-PER'),
 ('ca ngi', 'V', 'B-VP', 'O'),
 ('nhng', 'L', 'B-NP', 'O'),
 ('thành tu', 'N', 'B-NP', 'O'),
 ('vĩ di', 'A', 'B-AP', 'O'),
 ('ca', 'E', 'B-PP', 'O'),
 ('Liên Xô', 'Np', 'B-NP', 'B-LOC')]
```


CHAPTER 8

classify

`underthesea.classification.classify(X, domain=None)`

Text classification

Parameters

- `X({unicode, str})` – raw sentence
- `domain({None, 'bank'})` –
domain of text
 - None: general domain
 - bank: bank domain

Returns `tokens` – categories of sentence

Return type list

Examples

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import classify
>>> sentence = "HLV ngoi doi gn t mi thang dn dt tuyn Vit Nam"
>>> classify(sentence)
['The thao']
```

```
>>> sentence = "Tôi rt thích cách phc v ca nhân viên BIDV"
>>> classify(sentence, domain='bank')
('CUSTOMER SUPPORT',)
```


CHAPTER 9

sentiment

`underthesea.sentiment.sentiment (X, domain=None)`
Sentiment Analysis

Parameters

- `X ({unicode, str})` – raw sentence
- `domain ({None, 'bank'})` –
domain of text
 - None: general domain
 - bank: bank domain

Returns `tokens` – sentiment of sentence

Return type list

Examples

```
>>> # -*- coding: utf-8 -*-
>>> from underthesea import sentiment
>>> sentence = "Va smartbidv, va bidv online mà li k dùng chung 1 tài khon đăngnhp, rc ri!"
>>> sentiment(sentence, domain='bank')
('INTERNET BANKING#NEGATIVE',)
```


CHAPTER 10

Indices and tables

- genindex
- modindex
- search

Python Module Index

U

`underthesea.chunking`, 15

`underthesea.ner`, 17

`underthesea.pos_tag`, 13

`underthesea.sentiment`, 21

`underthesea.word_sent`, 11

Index

C

chunk() (in module underthesea.chunking), 15
classify() (in module underthesea.classification), 19

N

ner() (in module underthesea.ner), 17

P

pos_tag() (in module underthesea.pos_tag), 13

S

sentiment() (in module underthesea.sentiment), 21

U

underthesea.chunking (module), 15
underthesea.ner (module), 17
underthesea.pos_tag (module), 13
underthesea.sentiment (module), 21
underthesea.word_sent (module), 11

W

word_sent() (in module underthesea.word_sent), 11