

---

# **iosr-crawler Documentation**

***Release 1.0.0***

**Grzegorz Miejski, Mateusz Radko, Krzysztof Trzepla**

June 08, 2015



<b>1 Crawler Engine</b>	<b>3</b>
1.1 Crawler Engine . . . . .	3
1.2 DB Engine . . . . .	4
1.3 Search Engine . . . . .	5
<b>2 Extractor</b>	<b>7</b>
2.1 Extractor . . . . .	7
<b>3 User Interface</b>	<b>9</b>
3.1 Forms . . . . .	9
3.2 Models . . . . .	10
<b>4 Indices and tables</b>	<b>11</b>
<b>Python Module Index</b>	<b>13</b>



Contents:



---

## Crawler Engine

---

Contents:

### 1.1 Crawler Engine

```
class engine.CrawlerEngine.CrawlerEngine

class CustomSpider(*a, **kw)

    allowed_domains = ['en.wikipedia.org']
    config = {'start_urls': 'http://en.wikipedia.org/wiki/Programming_language', 'allowed_domains': 'en.wikipedia.org'}
    config_file = <closed file '/home/docs/checkouts/readthedocs.org/user_builds/iosr-crawler/checkouts/latest/src/engine/config.py' at 0x7fb1ab416450>
    config_path = '/home/docs/checkouts/readthedocs.org/user_builds/iosr-crawler/checkouts/latest/src/engine/config.py'
    crawler
    handles_request(request)
    log(message, level=10, **kw)
        Log the given messages at the given log level. Always use this method to send log messages from
        your spider
    make_requests_from_url(url)
    name = 'spider'
    parse(response)
    static parse_page(response)
    parse_start_url(response)
    process_results(response, results)
    rules = (<scrapy.contrib.spiders.crawl.Rule object at 0x7fb1ab416450>,)
    set_crawler(crawler)
    settings
    start_requests()
    start_urls = ['http://en.wikipedia.org/wiki/Programming_language']
```

CrawlerEngine.**add\_query** (*user\_id*, *query*)

Add crawling query for given user.

**Parameters**

- **user\_id** (*int*) – ID of user associated with the query.
- **query** (*str*) – User's query.

CrawlerEngine.**get\_urls** (*query*)

Retrieves all URLs associated with given query from database.

**Returns** list of URLs.

CrawlerEngine.**get\_user\_queries** (*user\_id*)

Retrieves user queries from database.

**Parameters** **user\_id** (*int*) – Id of user associated with the query.

**Returns** list of user queries.

**static** CrawlerEngine.**notify\_agents** ()

Notifies agent about new crawling query.

CrawlerEngine.**start\_crawling** ()

Notifies all agents and if crawling process is not started, starts it.

## 1.2 DB Engine

**class** engine.db\_engine.DbEngine

**add\_keywords** (*query*, *keywords*, *bucket\_name='keywords'*)

Adds keywords for given query to database.

**Parameters**

- **query** (*str*) – Query associated with keywords.
- **keywords** (*list*) – List of keywords produced from the query.

**add\_query** (*user\_id*, *query*, *bucket\_name='user\_queries'*)

Adds query to database.

**Parameters**

- **user\_id** (*int*) – Id of user associated with the query.
- **query** (*str*) – Query to be saved into database.

**add\_url** (*query*, *url*, *bucket\_name='urls'*)

Adds url for given query to database.

**Parameters**

- **query** (*str*) – Query associated with url.
- **url** (*str*) – URL of page satisfying search requirements.

**get\_all\_queries** (*bucket\_name='all\_queries'*)

Retrieves all queries from database.

**Returns** list of all queries.

**get\_keywords** (*query, bucket\_name='keywords'*)  
Retrieves all keywords associated with given query from database.

**Returns** list of keywords.

**get\_urls** (*query, bucket\_name='urls'*)  
Retrieves all URLs associated with given query from database.

**Returns** list of URLs.

**get\_user\_queries** (*user\_id, bucket\_name='user\_queries'*)  
Retrieves user queries from database.

**Parameters** **user\_id** (*int*) – Id of user associated with the query.

**Returns** list of user queries.

## 1.3 Search Engine

```
class engine.search_engine.SearchEngine

    reload_queries ()
        Reloads queries from database.

    search (content)
        Iterates over all queries and returns those for which number of found keywords satisfies search threshold.

        Parameters content (str) – content of web page associated with the URL.

        Returns list of queries for which search threshold was satisfied.

    search_in_url (url, content)
        Search web page content in order to find keywords.

        Parameters
            • url (str) – URL of web page being crawled.
            • content (str) – content of web page associated with the URL.
```



---

## Extractor

---

Contents:

### 2.1 Extractor

```
class nlp.extractor.NLPExtractor

    build_stop_word_regex()
        Creates stop word regex.

        Returns stop word pattern.

    static calculate_word_scores(phrase_list)
        Calculates words scores based on their frequency and degree.

        Parameters phrase_list (list) – List of phrases to be processed.

        Returns mapping between word and its score.

    static generate_candidate_keyword_scores(phrase_list, word_score)
        Generates scores for candidate keywords.

        Parameters
            • phrase_list (list) – List of phrases to be processed.
            • word_score (map) – Mapping between word and its score.

        Returns mapping between phrases and their scores.

    static generate_candidate_keywords(sentence_list, stopword_pattern)
        Generates list of keywords candidates.

        Parameters
            • sentence_list (list) – List of sentences to be processed.
            • stopword_pattern (str) – Stop words pattern.

        Returns list of keywords

    static is_number(word)
        Checks whether word is a number.

        Parameters word (str) – Word to be checked.

        Returns True or False
```

---

**load\_stop\_words ()**

Utility function to load stop words from a file and return as a list of words.

**Returns** list A list of stop words.

**run (text)**

Extracts keywords from the text.

**Parameters** **text** (*str*) – Text to be processed.

**Returns** list of keywords.

**static separate\_words (text, min\_word\_return\_size)**

Utility function to return a list of all words that are have a length greater than a specified number of characters.

**Parameters**

- **text** (*str*) – The text that must be split in to words.
- **min\_word\_return\_size** (*int*) – The minimum no of characters a word must have to be included.

**Returns** list of separated words.

**static split\_sentences (text)**

Utility function to return a list of sentences.

**Parameters** **text** (*str*) – The text that must be split in to sentences.

**Returns** sentences List of sentences created due to split.

---

## User Interface

---

Contents:

### 3.1 Forms

```
class ui.forms.QueryForm(data=None, files=None, auto_id=u'id_%s', prefix=None, initial=None,
                         error_class=<class 'django.forms.utils.ErrorList'>, label_suffix=None,
                         empty_permitted=False)
```

#### **add\_error**(*field*, *error*)

Update the content of *self.\_errors*.

The *field* argument is the name of the field to which the errors should be added. If its value is None the errors will be treated as NON\_FIELD\_ERRORS.

The *error* argument can be a single error, a list of errors, or a dictionary that maps field names to lists of errors. What we define as an “error” can be either a simple string or an instance of ValidationError with its message attribute set and what we define as list or dictionary can be an actual *list* or *dict* or an instance of ValidationError with its *error\_list* or *error\_dict* attribute set.

If *error* is a dictionary, the *field* argument *must* be None and errors will be added to the fields that correspond to the keys of the dictionary.

#### **add\_initial\_prefix**(*field\_name*)

Add a ‘initial’ prefix for checking dynamic initial values

#### **add\_prefix**(*field\_name*)

Returns the field name with a prefix appended, if this Form has a prefix set.

Subclasses may wish to override.

#### **as\_p**()

Returns this form rendered as HTML <p>s.

#### **as\_table**()

Returns this form rendered as HTML <tr>s – excluding the <table></table>.

#### **as\_ul**()

Returns this form rendered as HTML <li>s – excluding the <ul></ul>.

```
base_fields = OrderedDict([('query', <django.forms.fields.CharField object at 0x7fb1aaf4cb10>)])
```

#### **changed\_data**

**clean()**

Hook for doing any extra form-wide cleaning after Field.clean() has been called on every field. Any ValidationError raised by this method will not be associated with a particular field; it will have a special-case association with the field named ‘`__all__`’.

**declared\_fields = OrderedDict([('query', <django.forms.fields.CharField object at 0x7fb1aaf4cb10>)])****errors**

Returns an ErrorDict for the data provided for the form

**full\_clean()**

Cleans all of self.data and populates self.\_errors and self.cleaned\_data.

**has\_changed()**

Returns True if data differs from initial.

**has\_error (field, code=None)****hidden\_fields()**

Returns a list of all the BoundField objects that are hidden fields. Useful for manual form layout in templates.

**is\_multipart()**

Returns True if the form needs to be multipart-encoded, i.e. it has FileInput. Otherwise, False.

**is\_valid()**

Returns True if the form has no errors. Otherwise, False. If errors are being ignored, returns False.

**media****non\_field\_errors()**

Returns an ErrorList of errors that aren’t associated with a particular field – i.e., from Form.clean(). Returns an empty ErrorList if there are none.

**visible\_fields()**

Returns a list of BoundField objects that aren’t hidden fields. The opposite of the hidden\_fields() method.

## 3.2 Models

## **Indices and tables**

---

- genindex
- modindex
- search



**e**

engine.CrawlerEngine, 3  
engine.db\_engine.DbEngine, 4  
engine.search\_engine.SearchEngine, 5

**n**

nlp.extractor, 7

**u**

ui.forms, 9  
ui.models, 10



## A

add\_error() (ui.forms.QueryForm method), 9  
add\_initial\_prefix() (ui.forms.QueryForm method), 9  
add\_keywords() (engine.db\_engine.DbEngine.DbEngine method), 4  
add\_prefix() (ui.forms.QueryForm method), 9  
add\_query() (engine.CrawlerEngine.CrawlerEngine method), 3  
add\_query() (engine.db\_engine.DbEngine.DbEngine method), 4  
add\_url() (engine.db\_engine.DbEngine.DbEngine method), 4  
allowed\_domains (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
as\_p() (ui.forms.QueryForm method), 9  
as\_table() (ui.forms.QueryForm method), 9  
as\_ul() (ui.forms.QueryForm method), 9

## B

base\_fields (ui.forms.QueryForm attribute), 9  
build\_stop\_word\_regex() (nlp.extractor.NLPExtractor method), 7

## C

calculate\_word\_scores() (nlp.extractor.NLPExtractor static method), 7  
changed\_data (ui.forms.QueryForm attribute), 9  
clean() (ui.forms.QueryForm method), 9  
config (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
config\_file (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
config\_path (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
crawler (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
CrawlerEngine (class in engine.CrawlerEngine), 3  
CrawlerEngine.CustomSpider (class in engine.CrawlerEngine), 3

## D

DbEngine (class in engine.db\_engine.DbEngine), 4  
declared\_fields (ui.forms.QueryForm attribute), 10

## E

engine.CrawlerEngine (module), 3  
engine.db\_engine.DbEngine (module), 4  
engine.search\_engine.SearchEngine (module), 5  
errors (ui.forms.QueryForm attribute), 10

## F

full\_clean() (ui.forms.QueryForm method), 10

## G

generate\_candidate\_keyword\_scores()  
(nlp.extractor.NLPExtractor static method), 7  
generate\_candidate\_keywords()  
(nlp.extractor.NLPExtractor static method), 7  
get\_all\_queries() (engine.db\_engine.DbEngine.DbEngine method), 4  
get\_keywords() (engine.db\_engine.DbEngine.DbEngine method), 4

get\_urls() (engine.CrawlerEngine.CrawlerEngine method), 4  
get\_urls() (engine.db\_engine.DbEngine.DbEngine method), 5  
get\_user\_queries() (engine.CrawlerEngine.CrawlerEngine method), 4  
get\_user\_queries() (engine.CrawlerEngine.CrawlerEngine method), 4  
get\_user\_queries() (engine.db\_engine.DbEngine.DbEngine method), 5

## H

handles\_request() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3  
has\_changed() (ui.forms.QueryForm method), 10  
has\_error() (ui.forms.QueryForm method), 10  
hidden\_fields() (ui.forms.QueryForm method), 10

**I**

is\_multipart() (ui.forms.QueryForm method), 10  
is\_number() (nlp.extractor.NLPExtractor static method), 7  
is\_valid() (ui.forms.QueryForm method), 10

**L**

load\_stop\_words() (nlp.extractor.NLPExtractor method), 7  
log() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3

**M**

make\_requests\_from\_url() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3  
media (ui.forms.QueryForm attribute), 10

**N**

name (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
nlp.extractor (module), 7  
NLPExtractor (class in nlp.extractor), 7  
non\_field\_errors() (ui.forms.QueryForm method), 10  
notify\_agents() (engine.CrawlerEngine.CrawlerEngine static method), 4

**P**

parse() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3  
parse\_page() (engine.CrawlerEngine.CrawlerEngine.CustomSpider static method), 3  
parse\_start\_url() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3  
process\_results() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3

**Q**

QueryForm (class in ui.forms), 9

**R**

reload\_queries() (engine.search\_engine.SearchEngine.SearchEngine method), 5  
rules (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
run() (nlp.extractor.NLPExtractor method), 8

**S**

search() (engine.search\_engine.SearchEngine.SearchEngine method), 5  
search\_in\_url() (engine.search\_engine.SearchEngine.SearchEngine method), 5

**U**

SearchEngine (class in engine.search\_engine.SearchEngine), 5  
separate\_words() (nlp.extractor.NLPExtractor static method), 8  
set\_crawler() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3  
settings (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3  
split\_sentences() (nlp.extractor.NLPExtractor static method), 8  
start\_crawling() (engine.CrawlerEngine.CrawlerEngine method), 4  
start\_requests() (engine.CrawlerEngine.CrawlerEngine.CustomSpider method), 3  
start\_urls (engine.CrawlerEngine.CrawlerEngine.CustomSpider attribute), 3

**V**

visible\_fields() (ui.forms.QueryForm method), 10