

---

# **Nordlys Documentation**

*Release 0.2*

**IAI group - University of Stavanger**

**Mar 10, 2021**



---

## Contents

---

<b>1 Nordlys can be used ...</b>	<b>3</b>
<b>2 Contents</b>	<b>5</b>
<b>3 Indices and tables</b>	<b>71</b>
<b>Python Module Index</b>	<b>73</b>
<b>Index</b>	<b>75</b>



Nordlys is a toolkit for entity-oriented and semantic search, created by the IAI group at the University of Stavanger.

Entities (such as people, organizations, or products) are meaningful units for organizing information and can provide direct answers to many search queries. Nordlys is a toolkit for entity-oriented and semantic search.

Nordlys supports 3 functionalities in the context of entity-oriented search:

- Entity retrieval: Returns a ranked list of entities in response to a query
- Entity linking: Identifies entities in a query and links them to the corresponding entry in the Knowledge base
- Target type identification: Detects the target types (or categories) of a query

Check our [Web interface documentation](#) for illustration of each of these functionalities.



# CHAPTER 1

---

Nordlys can be used ...

---

- through a *web-based GUI*
- through a *RESTful API*
- as a *command line tool*
- as a *Python package*



## 2.1 Installation

Nordlys is a general-purpose semantic search toolkit, which can be used either as a Python package, as a command line tool, or as a service.

Data are a first-class citizen in Nordlys. To make use of the full functionality, the required data backend (MongoDB and Elasticsearch) need to be set up and data collections need to be loaded into them. There is built-in support for specific data collections, including DBpedia and Freebase. You may use the data dumps we prepared, or download, process and index these datasets from the raw sources.

### 2.1.1 1 Installing the Nordlys package

This step is required for all usages of Nordlys (i.e., either as a Python package, command line tool or service).

#### 1.1 Environment

Nordlys requires Python 3.5+ and a Python environment you can install packages in. We highly recommend using an [Anaconda Python distribution](#).

#### 1.2 Obtain source code

You can clone the Nordlys repo using the following:

```
$ git clone https://github.com/iai-group/nordlys.git
```

#### 1.3 Install prerequisites

Install Nordlys prerequisites using pip:

```
$ pip install -r requirements.txt
```

If you don't have pip yet, install it using

```
$ easy_install pip
```

---

**Note:** On Ubuntu, you might need to install lxml using a package manager

```
$ apt-get install python-lxml
```

---

### 1.4 Test installation

You may check if your installation has been successful by running any of the `cmd_usage`, e.g., from the root of your Nordlys folder issue

```
$ python -m nordlys.core.retrieval.retrieval
```

Alternatively, you can try importing Nordlys into a Python project. Make sure your local Nordlys copy is on the `PYTHONPATH`. Then, you may try, e.g., from `nordlys.core.retrieval.scorer` import `Scorer`.

Mind that this step is only to check the Python dependencies. It is rather limited what you can do with Nordlys without setting up data backed and loading the data components.

### 2.1.2 2 Setting up data backend

We use MongoDB and Elasticsearch for storing and indexing data. You can either connect to these services already running on some server or set these up on your local machine.

#### 2.1 MongoDB

If you need to install MongoDB yourself, follow the [instructions here](#).

Adjust the settings in `config/mongo.json`, if needed.

If you're using macOS, you'll likely need to change the soft limit of maxfiles to at least 64000, for mongoDB to work properly. Check the maxfiles limit of your system using:

```
$ launchctl limit
```

#### 2.2 Elasticsearch

If you need to install Elasticsearch yourself, follow the [instructions here](#). Note that Elasticsearch requires Java version 8.

Adjust the settings in `config/elastic.json`, if needed.

Nordlys has been tested with Elasticsearch version 5.x (and would likely need updating for newer version).

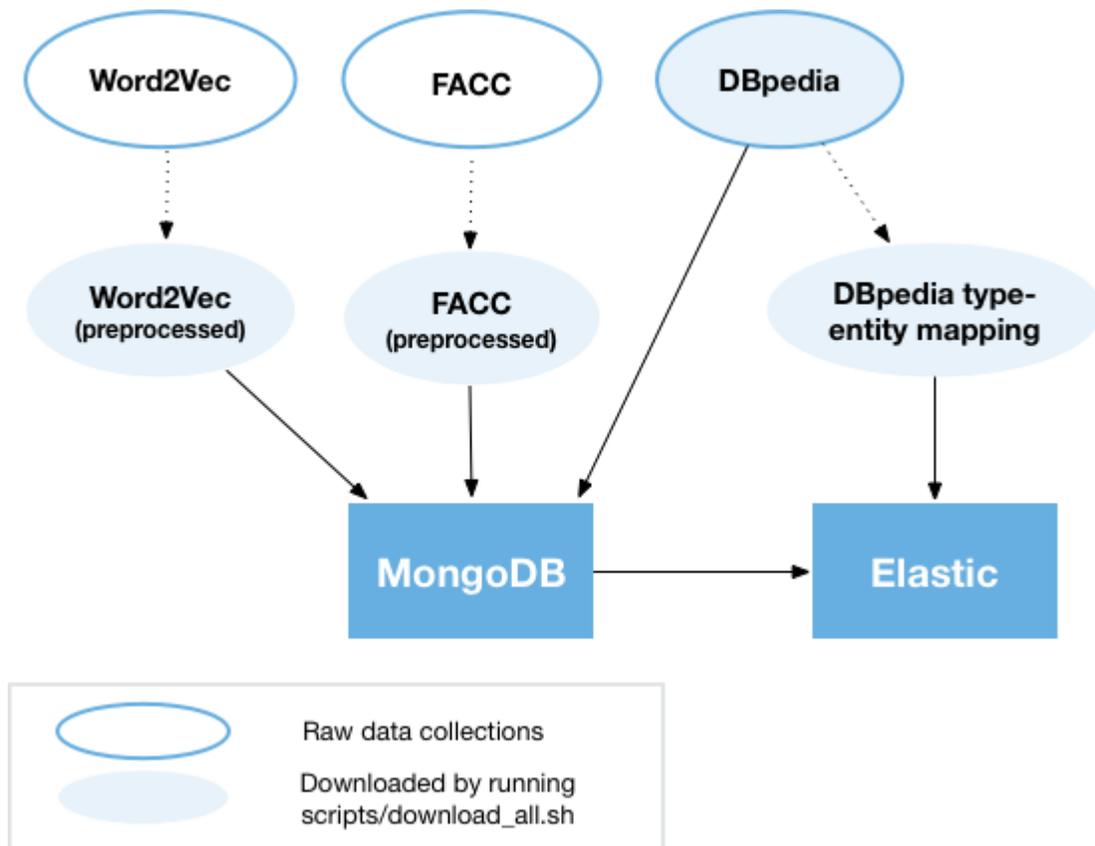
### 2.1.3 3 Loading data components

Data are a crucial component of Nordlys. While most of the functionality is agnostic of the underlying knowledge base, there is built-in support for working with specific data sources. This primarily means DBpedia, with associated resources from Freebase.

Specifically, Nordlys is shipped with functionality designed around **DBpedia 2015-10**, and the dumps we provide are for that particular version. However, we provide config files and instructions for working with **DBpedia 2016-10**, should you prefer a newer version.

Note that you may need only a certain subset of the data, depending on the required functionality. See *Data components* for a detailed description.

The figure below shows an overview of data sources and their dependencies.



#### 3.1 Load data to MongoDB

You can either load the data to MongoDB (i) from dumps that we made available or (ii) from the raw source files (DBpedia, FACC, Word2vec, etc.). Below, we discuss the former option. For the latter, see *Building MongoDB sources from raw data*. Note that processing from the raw sources takes significantly longer because of the nontrivial amount of data.

To load the data to MongoDB, you need to run the following commands from the main Nordlys folder. Note that the first dump is required for the core Nordlys functionality over DBpedia. The other dumps are optional, depending on whether the respective functionality is needed.

Command	Required for
<code>./scripts/load_mongo_dumps.sh mongo_dbpedia-2015-10.tar.bz2</code>	All
<code>./scripts/load_mongo_dumps.sh mongo_surface_forms_dbpedia.tar.bz2</code> <code>./scripts/load_mongo_dumps.sh mongo_surface_forms_facc.tar.bz2</code> <code>./scripts/load_mongo_dumps.sh mongo_fb2dbp-2015-10.tar.bz2</code>	EL and EC
<code>./scripts/load_mongo_dumps.sh mongo_word2vec-googlenews.tar.bz2</code>	TTI

### 3.2 Download auxiliary data files

The following files are needed for various services. You may download them all using

```
$ ./scripts/download_auxiliary.sh
```

Description	Location (relative to main Nordlys folder)	Required for
Type-to-entity mapping	<code>data/raw-data/dbpedia-2015-10/type2entity-mapping</code>	TTI
Freebase-to-DBpedia mapping	<code>data/raw-data/dbpedia-2015-10/freebase2dbpedia</code>	EL
Entity snapshot	<code>data/el</code>	EL <sup>1</sup>

- <sup>1</sup> If entity annotations are to be limited to a specific set; this file contains the proper named entities in DBpedia 2015-10

### 3.3 Build Elastic indices

There are multiple `elastic_indices` created for supporting different services. Run the following commands from the main Nordlys folder to build the indices for the respective functionality.

Command	Source	Required for
<code>./scripts/build_dbpedia_index.sh core</code>	MongoDB	ER, EL, TTI
<code>./scripts/build_dbpedia_index.sh types</code>	Raw DBpedia files <sup>1</sup>	TTI
<code>./scripts/build_dbpedia_index.sh uri</code>	MongoDB	ER <sup>2</sup>

- <sup>1</sup> Requires short entity abstracts and instance types files
- <sup>2</sup> only for ELR model

---

**Note:** To use the 2016-10 version of DBpedia, add `2016-10` as a 2nd argument to the above scripts.

---

## 2.2 Data components

### 2.2.1 Data sources

## DBpedia

We use DBpedia as the main underlying knowledge base. In particular, we prepared dumps for DBpedia version 2015-10.

DBpedia is distributed, among other formats, as a set of .ttl.bz2 files. We use a selection of these files, as defined in `data/config/dbpedia2mongo.config.json`. You can download these files from [DBpedia Website](#) directly or running `./scripts/download_dbpedia.sh` from the main Nordlys folder. Running the script will place the downloaded files under `data/raw-data/dbpedia-2015-10/`.

We also provide a minimal sample from DBpedia under `data/dbpedia-2015-10-sample`, which can be used for testing/development in a local environment.

## FACC

The Freebase Annotations of the ClueWeb Corpora (FACC) is used for building entity surface form dictionary. You can download the collection from its [main Website](#). and further process it using our scripts. Alternatively, you can download the preprocessed data from our server. Check the README file under `data/raw-data/facc` for detailed information.

## Word2Vec

Word2Vec vectors (300D) trained on Google News corpus, which can be downloaded from the its [Website](#). Check the README file under `data/raw-data/word2vec` for detailed information.

### 2.2.2 MongoDB collections

The table below provides an overview of the MongoDB collections that are used by the different services.

Name	Description	EC	ER	EL	TTI
<code>dbpedia-2015-10</code>	DBpedia	+ <sup>1</sup>	+ <sup>2</sup>		+ <sup>3</sup>
<code>fb2dbp-2015-10</code>	Mapping from Freebase to DBpedia IDs	+ <sup>4</sup>		+	
<code>surface_forms_dbpedia</code>	Entity surface forms from DBpedia	+ <sup>5</sup>		+ <sup>6</sup>	
<code>surface_forms_facc</code>	Entity surface forms from FACC	+ <sup>7</sup>		+	
<code>word2vec-googlenews</code>	Word2vec trained on Google News				+ <sup>8</sup>

- <sup>1</sup> for entity ID-based lookup and DBpedia2Freebase mapping functionalities
- <sup>2</sup> only for building the Elastic entity index; not used later in the retrieval process
- <sup>3</sup> for entity-centric TTI method
- <sup>4</sup> for Freebase2DBpedia mapping functionality
- <sup>5</sup> for entity surface form lookup from DBpedia
- <sup>6</sup> for all EL methods other than “commonness”
- <sup>7</sup> for entity surface form lookup from FACC
- <sup>8</sup> for LTR TTI method

## Building MongoDB sources from raw data

To build the above tables from raw data (as opposed to the provided dumps), first make sure that you have the raw data files.

- For DBpedia, these may be downloaded using `./scripts/download_dbpedia.sh`
- For the FACC and Word2vec data files, execute `./scripts/download_raw.sh`

To load DBpedia to MongoDB, run

```
python -m nordlys.core.data.dbpedia.dbpedia2mongo data/config/dbpedia-2015-10/  
↳dbpedia2mongo.config.json
```

---

**Note:** To use the DBpedia 2015-10 sample shipped with Nordlys, as opposed to the full collection, change the value path to `data/raw-data/dbpedia-2015-10_sample/` in `dbpedia2mongo.config.json`.

---

### 2.2.3 Elastic indices

Name	Description	ER	EL	TTI
<code>dbpedia_2015_10</code>	DBpedia index	+	+ <sup>1</sup>	+ <sup>2</sup>
<code>dbpedia_2015_10_uri</code>	DBpedia URI-only index	+ <sup>3</sup>		
<code>dbpedia_2015_10_types</code>	DBpedia types index			+ <sup>4</sup>

- <sup>1</sup> for all EL methods other than “commonness”
- <sup>2</sup> only for entity-centric TTI method
- <sup>3</sup> only for ELR entity ranking method
- <sup>4</sup> only for type-centric TTI method

## 2.3 RESTful API

The following Nordlys services can be accessed through a RESTful API:

- *Entity Retrieval*
- *Entity Linking in Queries*
- *Target Type Identification*
- *Entity Catalog*

Below we describe the usage for each of the services.

The Nordlys endpoint URL is <http://api.nordlys.cc/>.

### 2.3.1 Entity Retrieval

The service presents a ranked list of entities in response to an entity-bearing query.

## Endpoint URI

```
http://api.nordlys.cc/er
```

## Example

### Request:

```
http://api.nordlys.cc/er?q=total+recall&1st_num_docs=100&model=lm
```

### Response:

```
{
  "query": "total recall",
  "total_hits": 1000,
  "results": {
    "0": {
      "entity": "<dbpedia:Total_Recall_(1990_film)>",
      "score": -10.042525028471253
    },
    "1": {
      "entity": "<dbpedia:Total_Recall_(2012_film)>",
      "score": -10.295316626850521
    },
    ...
  }
}
```

## Parameters

The following table lists the parameters needed in the request URL for entity retrieval.

Parameters	
q ( <i>required</i> )	Search query
1st_num_docs	The number of documents that will be re-ranked using a model. The recommended value (esp. for baseline comparisons) is 1000. Lower values, like 100, are recommended only when efficiency matters ( <i>default: 1000</i> ).
start	Starting offset for ranked documents ( <i>default: 0</i> ).
fields_return	Comma-separated list of fields to return for each hit ( <i>default: ""</i> ).
model	Name of the retrieval model; Accepted values: "bm25", "lm", "mlmprms" ( <i>default: "lm"</i> ). <ul style="list-style-type: none"> <li>• BM25: The BM25 model, as implemented in Elasticsearch. This is the most efficient that is provided by Nordlys (to this date).</li> <li>• LM: Language Modeling [1] approach, which employs a single a single field representation of entities.</li> <li>• MLM: The Mixture of Language Models [2], which uses a linear combination of language models built for each field.</li> <li>• PRMS: The Probabilistic Model for Semistructured Data [3], which uses collection statistics to compute field weights in for MLM model.</li> </ul>
field	The name of the field used for LM ( <i>default: "catchall"</i> ).
fields	Comma-separated list of the fields for PRMS ( <i>default: "catchall"</i> ).
field_weights	Comma-separated list of fields and their corresponding weights for MLM ( <i>default: "catchall:1"</i> ).
smoothing_method	Smoothing method for LM-based models; Accepted values: jm, dirichlet ( <i>default: "dirichlet"</i> ).
smoothing_param	The value of smoothing parameters (lambda or mu); Accepted values: float, "avg_len" ( <i>default for "jm": 0.1, default for "dirichlet": 2000</i> ).

### 2.3.2 Entity Linking in Queries

The service identifies entities in queries and links them to the corresponding entry in the Knowledge base (DBpedia).

#### Endpoint URI

```
http://api.nordlys.cc/el
```

#### Example

- **Request:** <http://api.nordlys.cc/el?q=total+recall>
- **Response:**

```

{
  "processed_query": "total recall",
  "query": "total recall",
  "results": [
    {
      "entity": "<dbpedia:Total_Recall_(1990_film)>",
      "mention": "total recall",
      "score": 0.4013333333333334
    },
    {
      "entity": "<dbpedia:Total_Recall_(2012_film)>",
      "mention": "total recall",
      "score": 0.315
    }
  ]
}

```

## Parameters

The following table lists the parameters needed in the request URL for entity linking.

Parameters	
q ( <i>required</i> )	The search query
method	The name of method; Accepted values ( <i>default</i> : "cmns") <ul style="list-style-type: none"> <li>cmns: The baseline method that uses the overall popularity of entities as link targets, implemented based on [5].</li> <li>ltr: The learning-to-rank model, implemented based on the LTR-greedy in [9]. Note that the implemented method is slightly different from [9] (due to efficiency reasons).</li> </ul>
threshold	The entity linking threshold ( <i>default</i> : 0.1).

### 2.3.3 Target Type Identification

The service assigns target types (or categories) to queries from the DBpedia type taxonomy.

#### Endpoint URI

```
http://api.nordlys.cc/tti
```

#### Example

- **Request:** <http://api.nordlys.cc/tti?q=obama>
- **Response:**

```

{
  "query": "obama",
  "results": {
    "0": {
      "score": 3.3290777,
      "type": "<dbo:Ambassador>"
    },
    "1": {
      "score": 3.2955842,
      "type": "<dbo:Election>"
    },
    ...
  }
}

```

## Parameters

The following table lists the parameters needed in the request URL for target type identification.

Parameters	
q ( <i>required</i> )	The search query
method	The name of method; accepted values: “tc”, “ec”, “ltr” ( <i>default</i> : “tc”). <ul style="list-style-type: none"> <li>• TC: The Type Centric (TC) method based on [6]. Both BM25 and LM models can be used as a retrieval model here. This method fits the early fusion design pattern in [7].</li> <li>• EC: The Entity Centric (EC) method, as described in [6]. Both BM25 and LM models can be used as a retrieval model here. This method fits the late fusion design pattern in [7].</li> <li>• LTR: The Learning-To-Rank (LTR) method, as proposed in [8].</li> </ul>
num_docs	The number of top ranked target types to retrieve ( <i>default</i> : 10).
start	The starting offset for ranked types.
model	Retrieval model, if method is “tc” or “ec”; Accepted values: “lm”, “bm25”.
ec_cutoff	If method is “ec”, rank cut-off of top-K entities for EC TTL.
field	Field name, if method is “tc” or “ec”.
smoothing_method	If model is “lm”, smoothing method; accepted values: “jm”, “dirichlet”.
Smoothing_param	If model is “lm”, smoothing parameter; accepted values: float, “avg_len”.

### 2.3.4 Entity Catalog

This service is used for representing entities (with IDs, name variants, attributes, and relationships). Additionally, it provides statistics that can be utilized, among others, for result presentation (e.g., identifying prominent properties when generating entity cards).

## Endpoint URI

```
http://api.nordlys.cc/ec
```

## Look up entity by ID

- **Request:**

```
http://api.nordlys.cc/ec/lookup_id/<entity_id>
```

- **Example:**

```
http://api.nordlys.cc/ec/lookup_id/<dbpedia:Albert_Einstein>
```

- **Response:**

```
{
  "<dbo:abstract>": ["Albert Einstein was a German-born theoretical physicist ..
→. ],
  "<dbo:academicAdvisor>": ["<dbpedia:Heinrich_Friedrich_Weber>"],
  "<dbo:almaMater>": [
    "<dbpedia:ETH_Zurich>",
    "<dbpedia:University_of_Zurich>"
  ],
  "<dbo:award>": [
    "<dbpedia:Nobel_Prize_in_Physics>",
    "<dbpedia:Max_Planck_Medal>",
    ...
  ],
  "<dbo:birthDate>": ["1879-03-14"],
  ...
}
```

## Look up entity by name (DBpedia)

Looks up an entity by its surface form in DBpedia.

- **Request:**

```
http://api.nordlys.cc/ec/lookup_sf/dbpedia/<sf>
```

- **Example:**

```
http://api.nordlys.cc/ec/lookup_sf/dbpedia/new%20york
```

- **Response:**

```
{
  "_id": "new york"
  "<rdfs:label>" : {
    "<dbpedia:New_York>": 1
  }
  "<dbo:wikiPageDisambiguates>": {
    "<dbpedia:Manhattan>": 1,
    "<dbpedia:New_York,_Kentucky>": 1,
    ...
  }
}
```

(continues on next page)

(continued from previous page)

```
}  
  ...  
}
```

### Look up entity by name (FACC)

Looks up an entity by its surface form in FACC.

- **Request:**

```
http://api.nordlys.cc/ec/lookup_sf/facc/<sf>
```

- **Example:**

```
http://api.nordlys.cc/ec/lookup_sf/facc/new%20york
```

- **Response:**

```
{  
  "_id" : "new york",  
  "facc12" : {  
    "<fb:m.02_286>": 18706787,  
    "<fb:m.02_53fb>": 49,  
    "<fb:m.02_b91>": 87,  
    "<fb:m.02_143>": 12,  
    "<fb:m.02_15n>": 23963,  
    ...  
  }  
}
```

### Map Freebase entity ID to DBpedia ID

- **Request:**

```
http://api.nordlys.cc/ec/freebase2dbpedia/<fb_id>
```

- **Example:**

```
http://api.nordlys.cc/ec/freebase2dbpedia/<fb:m.02_286>
```

- **Response:**

```
{  
  "dbpedia_ids": [  
    "<dbpedia:New_York_City>"  
  ]  
}
```

### Map DBpedia entity ID to Freebase ID

- **Request:**

```
http://api.nordlys.cc/ec/dbpedia2freebase/<dbp_id>
```

- **Example:**

```
http://api.nordlys.cc/ec/dbpedia2freebase/<dbpedia:New_York>
```

- **Response:**

```
{
  "freebase_ids": [
    "<fb:m.059rby>"
  ]
}
```

## Parameters

The following table lists the parameters needed in the request URL for entity catalog.

Parameter	
entity id	It is in the form of “<dbpedia:XXX>”, where XXX denotes the DBpedia/Wikipedia ID of an entity.
sf	Entity surface form (e.g., “john smith”, “new york”). It needs to be url-escaped.
fb_id	Freebase ID
bdp_id	DBpedia ID

## 2.3.5 References

- [1] Jay M Ponte and W Bruce Croft . 1998. *A Language modeling approach to information retrieval*. In Proc. of SIGIR '98. 275–281.
- [2] Paul Ogilvie and Jamie Callan. 2003. *Combining document representations for known-item search*. Proc. of SIGIR '03 (2003), 143–150.
- [3] Jinyoung Kim, Xiaobing Xue, and W Bruce Croft . 2009. *A probabilistic retrieval model for semistructured data*. In Proc. of ECIR '09. 228–239.
- [4] Faegheh Hasibi, Krisztian Balog, and Svein Erik Bratsberg. 2016. *Exploiting entity linking in queries for entity retrieval*. In Proc. of ICTIR '16. 171–180. [\[BIB\]](#) [\[PDF\]](#)
- [5] Faegheh Hasibi, Krisztian Balog, and Svein Erik Bratsberg. 2015. *Entity linking in queries: Tasks and Evaluation*. In Proc. of ICTIR '15. 171–180. [\[BIB\]](#) [\[PDF\]](#)
- [6] Krisztian Balog, Robert Neumayer. 2012. *Hierarchical target type identification for entity-oriented queries*. In Proc. of CIKM '12. 2391–2394. [\[BIB\]](#) [\[PDF\]](#)
- [7] Shuo Zhang and Krisztian Balog. *Design Patterns for Fusion-Based Object Retrieval*. In Proc. of ECIR '17. 684-690. [\[BIB\]](#) [\[PDF\]](#)
- [8] Darío Garigliotti, Faegheh Hasibi, and Krisztian Balog. *Target Type Identification for Entity-Bearing Queries*. In Proc. of SIGIR '17. [\[BIB\]](#) [\[PDF\]](#)
- [9] Faegheh Hasibi, Krisztian Balog, and Svein Erik Bratsberg. 2017. *Entity linking in queries: Efficiency vs. Effectiveness*. In Proc. of ECIR '17. 40-53. [\[BIB\]](#) [\[PDF\]](#)

## 2.4 Web-based GUI

Nordlys is shipped with a web-based graphical user interface, which is built on the [Nordlys API](#). It is a wrapper for all functionalities provided by Nordlys toolkit and can be used, e.g., to perform user studies on result presentation.

The implementation of the Web interface is based on Flask and Bootstrap. Below we describe the functionalities provided by the excerpts from the Web GUI with their interface.

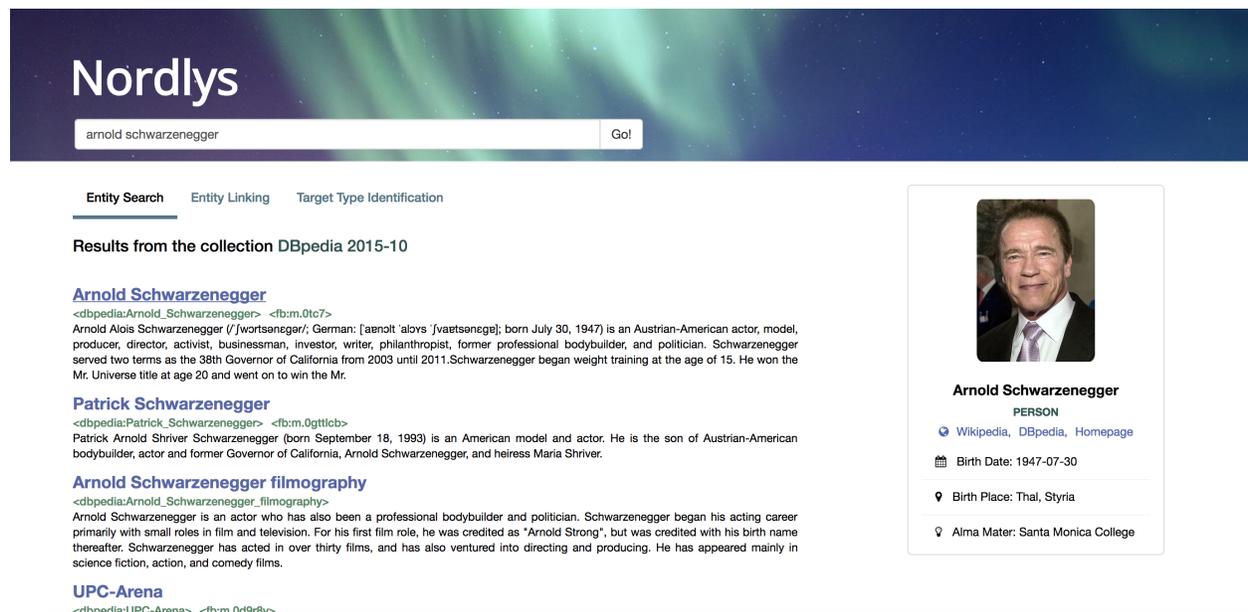
**Note:** The interface can be accessed via: <http://gui.nordlys.cc/>

### 2.4.1 Entity search

The entity search tab provides a ranked list of entities in response to an entity-bearing query. We generate the results by calling the Entity retrieval service of our API; e.g.,

```
http://api.nordlys.cc/er?q=total+recall&model=lm&1st_num_docs=100&fields_  
↪return=abstract
```

By clicking on each result, we present an entity card, containing the factual summary of the entity and an image. We use the entity catalog (EC) service to generate these cards.



The screenshot shows the Nordlys web interface. At the top, there is a search bar with the text "arnold schwarzenegger" and a "Go!" button. Below the search bar, there are three tabs: "Entity Search" (which is selected), "Entity Linking", and "Target Type Identification". The main content area displays "Results from the collection DBpedia 2015-10". There are three search results listed:

- Arnold Schwarzenegger**: A brief biography mentioning his birth date (July 30, 1947) and his career as an actor, model, producer, director, activist, businessman, investor, writer, philanthropist, former professional bodybuilder, and politician. It also notes he served two terms as the 38th Governor of California from 2003 until 2011.
- Patrick Schwarzenegger**: A brief biography mentioning his birth date (September 18, 1993) and his career as a model and actor. It notes he is the son of Arnold Schwarzenegger and heiress Maria Shriver.
- Arnold Schwarzenegger filmography**: A brief overview of his acting career, noting he has acted in over thirty films and has also ventured into directing and producing.

Below the search results, there is a section for "UPC-Arena" with a link to "UPC-Arena".

On the right side of the interface, there is an entity card for "Arnold Schwarzenegger". The card includes a portrait photo, the name "Arnold Schwarzenegger", the type "PERSON", and several links: "Wikipedia", "DBpedia", and "Homepage". It also lists "Birth Date: 1947-07-30", "Birth Place: Thal, Styria", and "Alma Mater: Santa Monica College".

Fig. 1: Nordlys Web interface - Entity Search

### 2.4.2 Entity linking in queries

For entity linking in queries, we use the baseline CMNS method, with threshold 0.1. For example, we make the following call to the entity linking service of our API:

```
http://api.nordlys.cc/el?q=arnold+schwarzenegger+total+recall&method=cmns&threshold=0.  
↪1
```

### 2.4.3 Target Type Identification

For target type identification, we employ the type centric method; e.g.,

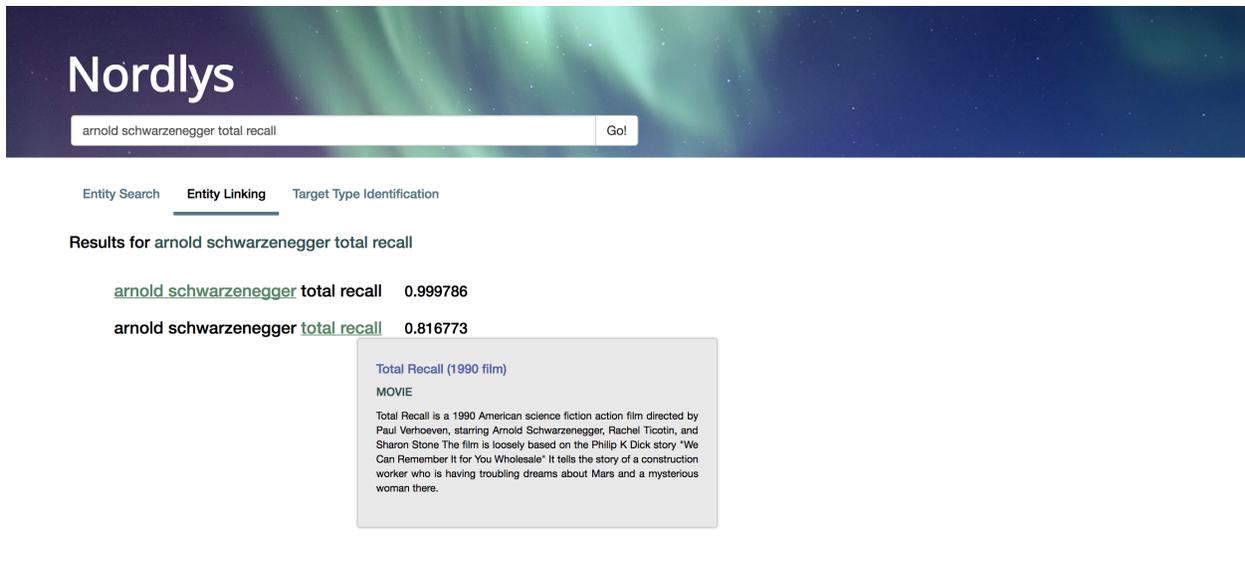


Fig. 2: Nordlys Web interface - Entity Linking

```
http://api.nordlys.cc/tti?q=obama
```

**Note:** For detailed information about our API calls, see `api_usage`

## 2.5 Command line usage

Nordlys provides the following command line applications:

- **Core applications**
  - Retrieval *general usage* and *configurations*
  - Machine learning *general usage* and configurations
- **Entity-oriented applications**
  - *Entity catalog*
  - *Entity retrieval*
  - Entity linking
  - *Target type identification*

## 2.6 Nordlys architecture

Nordlys is based on a *multitier architecture* with three layers:

- *core* (data tier)
- *logic*



Fig. 3: Nordlys Web interface - Target Type Identification

- *services* (presentation tier)

## 2.6.1 Core tier

The *core* tier provides basic functionalities and is connected to various third-party tools. The functionalities include:

- *Retrieval* (based on Elasticsearch)
- *Storage* (based on MongoDB)
- *Machine learning* (based on scikit-learn)
- *Evaluation* (based on trec-eval)

Additionally, a separate *data* package is provided with functionality for loading and preprocessing standard data sets (DBpedia, Freebase, ClueWeb, etc.).

It is possible to connect additional external tools (or replace our default choices) by implementing standard interfaces of the respective core modules.

---

**Note:** The core layer represents a versatile general-purpose modern IR library, which may also be accessed using command line tools.

---

## 2.6.2 Logic tier

The *logic* tier contains the main business logic, which is organized around five main modules:

1. *Entity* provides access to the entity catalog (including knowledge bases and entity surface form dictionaries).
2. *Query* provides the representation of search queries along with various preprocessing methods.
3. *Features* is a collection of entity-related features, which may be used across different search tasks.

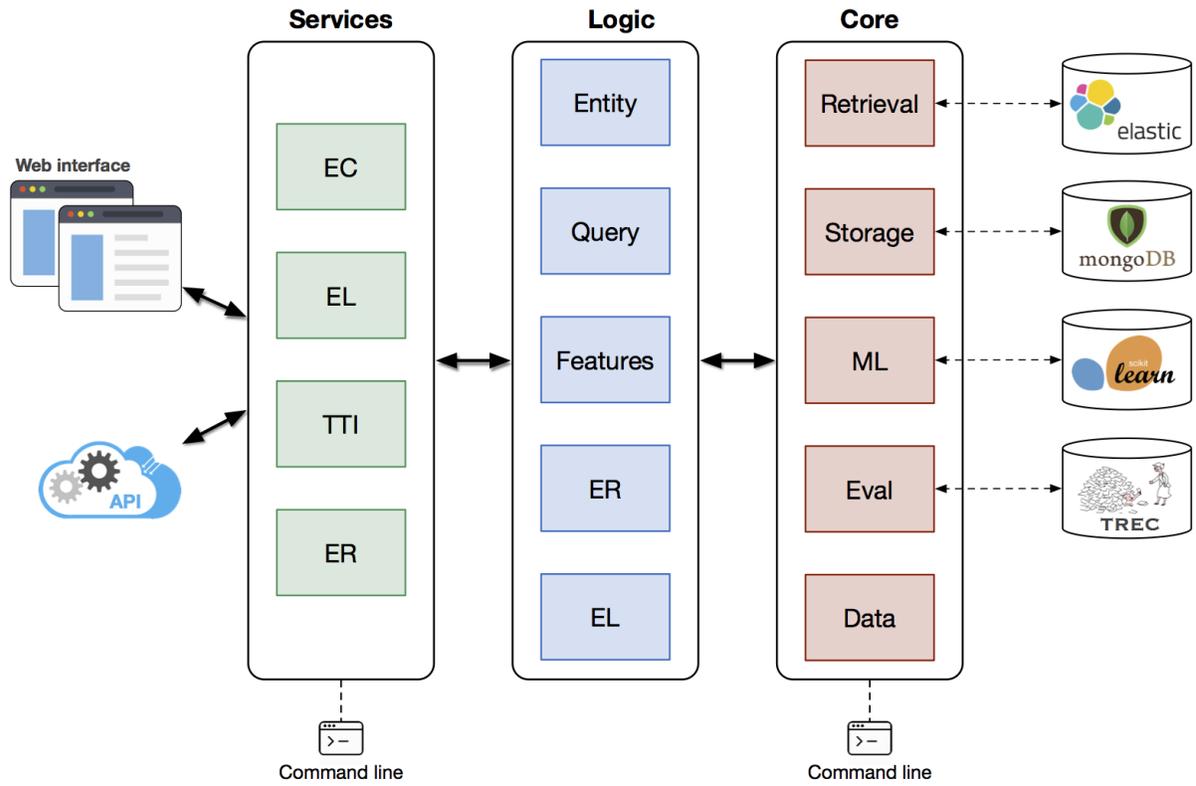


Fig. 4: Nordlys architecture.

4. *Entity retrieval* contains various entity ranking methods.
5. *Entity linking* implements entity linking functionality.

The logic layer may not be accessed directly (i.e., as a service or as a command line application).

### 2.6.3 Services tier

The *services* tier provides end-user access to the toolkit's functionality, throughout the *command line*, *API*, and *web interface*. Four main types of service is available:

1. *Entity retrieval*
2. *Entity linking*
3. *Target type identification*
4. *Entity catalog*

## 2.7 nordlys package

Nordlys toolkit is provided as the `nordlys` Python package.

Nordlys is delivered with an extensive and detailed documentation, from package-level overviews, to module usage examples, to fully documented classes and methods.

### 2.7.1 Subpackages

#### **nordlys.core package**

##### **Core packages**

These low-level packages (basic IR, ML, NLP, storage, etc.) are basic building blocks that do not have any interdependencies among each other.

##### **Subpackages**

#### **nordlys.core.data package**

##### **Subpackages**

#### **nordlys.core.data.dbpedia package**

##### **Submodules**

#### **nordlys.core.data.dbpedia.create\_sample module**

#### **nordlys.core.data.dbpedia.dbpedia2mongo module**

`nordlys.core.data.dbpedia.dbpedia_surfaceforms2mongo` module

`nordlys.core.data.dbpedia.freebase2dbpedia2mongo` module

`nordlys.core.data.dbpedia.indexer_dbpedia` module

`nordlys.core.data.dbpedia.indexer_dbpedia_types` module

## DBpedia Types Indexer

Builds a DBpedia type index from entity abstracts.

The index is build directly from DBpedia files in .ttl.bz2 format (i.e., MongoDB is not needed).

### Usage

```
python -m nordlys.core.dbpedia.indexer_dbpedia_types -c <config_file>
```

### Config parameters

- **index\_name**: name of the index
- **dbpedia\_files\_path**: path to DBpedia .ttl.bz2 files

**Authors** Krisztian Balog, Dario Garigliotti

**class** `nordlys.core.data.dbpedia.indexer_dbpedia_types.IndexerDBpediaTypes` (*config*)

Bases: `object`

**build\_index** (*force=False*)

Builds the index.

Note: since DBpedia only has a few hundred types, no bulk indexing is needed.

**Parameters** **force** – True iff it is required to overwrite the index (i.e. by creating it by force); False by default. :type force: bool :return:

**name**

`nordlys.core.data.dbpedia.indexer_dbpedia_types.arg_parser()`

`nordlys.core.data.dbpedia.indexer_dbpedia_types.main(args)`

`nordlys.core.data.dbpedia.indexer_dbpedia_uri` module

`nordlys.core.data.facc` package

### Submodules

`nordlys.core.data.facc.facc2mongo` module

### Facc to Mongo

Adds entity surface forms from the Freebase Annotated ClueWeb Corpora (FACC).

The input to this script is (name variant, Freebase entity, count) triples. See *data/facc1/README.md* for the preparation of FACC data in such format.

**Authors** Krisztian Balog, Faegheh Hasibi

```
class nordlys.core.data.facc.facc2mongo.FACCToMongo (config)
```

```
    Bases: object
```

```
    Inserts FACC surface forms to Mongo.
```

```
    build ()
```

```
        Builds surface form collection from FACC annotations.
```

```
nordlys.core.data.facc.facc2mongo.arg_parser ()
```

```
nordlys.core.data.facc.facc2mongo.main (args)
```

### nordlys.core.data.word2vec package

#### Submodules

#### nordlys.core.data.word2vec.word2vec2mongo module

### Word2vec to Mongo

Loads Word2Vec to MongoDB.

**Authors** Faegheh Hasibi, Dario Garigliotti

```
class nordlys.core.data.word2vec.word2vec2mongo.Word2VecToMongo (config)
```

```
    Bases: object
```

```
    build ()
```

```
        Builds word2vec collection from GoogleNews 300-dim pre-trained corpus.
```

```
nordlys.core.data.word2vec.word2vec2mongo.arg_parser ()
```

```
nordlys.core.data.word2vec.word2vec2mongo.main (args)
```

### nordlys.core.eval package

This is the (**intro** of the) eval package doc. It's in the `__init__.py` of the package, as an example for the other packages, and IT SHOULD BE cleaned up without all of this ;).

It could contain some/all eval hints.

It might contain, e.g., items in this way

- This
- That

and also code snippets like usual, using as usual these double colons and the indented code:

```
instance = Instance(id)
```

And that's it.

(In the following I included all the up-to-now available modules, sectioned by criterion at first sight.)

## Submodules

### nordlys.core.eval.eval module

#### Evaluation

Console application for eval package.

**Author** Faegheh Hasibi

```
class nordlys.core.eval.eval.Eval (operation, qrels=None, runs=None, metric=None, out-
                                     put_file=None)
```

Bases: `object`

Main entry point for eval package.

#### Parameters

- **operation** – operation (see `OPERATIONS` for allowed values)
- **qrels** – name of qrels file
- **runs** – name of run files
- **metric** – metric

```
OPERATIONS = ['query_diff']
```

```
OP_QUERY_DIFF = 'query_diff'
```

```
run ()
```

Runs the operation with the given arguments.

```
nordlys.core.eval.eval.arg_parser ()
```

```
nordlys.core.eval.eval.main (args)
```

### nordlys.core.eval.plot\_diff module

#### Plot Differences

Plots a series of scores which represent differences.

**Authors** Shuo Zhang, Krisztian Balog

```
class nordlys.core.eval.plot_diff.QueryDiff
```

Bases: `object`

```
SCORES = [25, 20, 10, 5, 0, -1, -5, -10]
```

```
create_pdf (diff_file, pdf_file, title="", xlabel="", ylabel="", aspect_ratio='equal', separator='\t')
```

Create bar plot for differences in pdf.

This function is used to load difference .csv file, create bar plot and store as a pdf file.

**Pdf** created and saved pdf file

**make\_plot** ()  
Make a bar plot using SCORES

## nordlys.core.eval.query\_diff module

### Query Differences

Computes query-level differences between two runs.

**Authors** Shuo Zhang, Krisztian Balog, Dario Garigliotti

**class** nordlys.core.eval.query\_diff.**QueryDiff** (*run1\_file, run2\_file, qrels, metric*)  
Bases: `object`

#### Parameters

- **run1\_file** – name of run1 file (baseline)
- **run2\_file** – name of run2 file (new method)
- **qrels** – name of qrels file
- **metric** – metric

#### Returns

**dump\_differences** (*output\_file*)  
Outputs query-level differences between two methods into a tab-separated file.

The first method is considered the baseline, the differences are with respect to that. Output format: queryID  
res1 res2 diff(res2-res1)

## nordlys.core.eval.trec\_eval module

### Trec Evaluation

Wrapper for trec\_eval.

**Authors** Dario Garigliotti, Shuo Zhang

**class** nordlys.core.eval.trec\_eval.**TrecEval**  
Bases: `object`

Holds evaluation results obtained using trec\_eval.

**evaluate** (*qrels\_file, run\_file, eval\_file=None*)  
Evaluates a runfile using trec\_eval. Optionally writes evaluation output to file.

#### Parameters

- **qrels\_file** – name of qrels file
- **run\_file** – name of run file
- **eval\_file** – name of evaluation output file

**get\_query\_ids** ()  
Returns the set of queryIDs for which we have results.

**get\_score** (*query\_id*, *metric*)

Returns the score for a given queryID and metric.

**Parameters**

- **query\_id** – queryID
- **metric** – metric

**Returns** score (or None if not found)

**load\_results** (*eval\_file*)

Loads results from an existing evaluation file.

**Parameters** **eval\_file** – name of evaluation file

## nordlys.core.eval.trec\_qrels module

### Trec Qrels

Utility module for working with TREC qrels files.

### Usage

**Get statistics about a qrels file** `trec_qrels <qrels_file> -o stat`

**Filter qrels to contain only documents from a given set** `trec_qrels <qrels_file> -o filter_docs -d <doc_ids_file> -f <output_file>`

**Filter qrels to contain only queries from a given set** `trec_qrels <qrels_file> -o filter_qs -q <query_ids_file> -f <output_file>`

**Author** Krisztian Balog

**class** `nordlys.core.eval.trec_qrels.TrecQrels` (*file\_name=None*)

Bases: `object`

Represents relevance judgments (TREC qrels).

**filter\_by\_doc\_ids** (*doc\_ids\_file*, *output\_file*)

Filters qrels for a set of selected docIDs and outputs the results to a file.

**Parameters**

- **doc\_ids\_file** – File with one docID per line
- **output\_file** – Output file name

**filter\_by\_query\_ids** (*query\_ids\_file*, *output\_file*)

Filters qrels for a set of selected queryIDs and outputs the results to a file.

**Parameters**

- **query\_ids\_file** – File with one queryID per line
- **output\_file** – Output file name

**get\_queries** ()

Returns the set of queries.

**get\_rel** (*query\_id*)

Returns relevance level for a given query.

**Parameters** `query_id` – queryID

**Returns** dict (docID as key and relevance as value) or None

**load** (*file\_name*)

Loads qrels from file.

**Parameters** `file_name` – name of qrels file

**num\_rel** (*query\_id*, *min\_rel=1*)

Returns the number of relevant results for a given query.

**Parameters**

- `query_id` – queryID
- `min_rel` – minimum relevance level

**Returns** number of relevant results

**print\_stat** ()

Prints simple statistics.

`nordlys.core.eval.trec_qrels.arg_parser()`

`nordlys.core.eval.trec_qrels.main(args)`

## nordlys.core.eval.trec\_run module

### Trec run

Utility module for working with TREC runfiles.

### Usage

**Get statistics about a runfile** `trec_run <run_file> -o stat`

**Filter runfile to contain only documents from a given set** `trec_run <run_file> -o filter -d <doc_ids_file> -f <output_file> -n <num_results>`

**Authors** Krisztian Balog, Dario Garigliotti

**class** `nordlys.core.eval.trec_run.TrecRun` (*file\_name=None*, *normalize=False*,  
*remap\_by\_exp=False*, *run\_id=None*)

Bases: `object`

Represents a TREC runfile.

**Parameters**

- `file_name` – name of the run file
- `normalize` – whether retrieval scores are to be normalized for each query (default: False)
- `remap_by_exp` – whether scores are to be converted from the log-domain by taking their exp (default: False)

**filter** (*doc\_ids\_file*, *output\_file*, *num\_results=100*)

Filters runfile to include only selected docIDs and outputs the results to a file.

**Parameters**

- `doc_ids_file` – file with one doc\_id per line

- **output\_file** – output file name
- **num\_results** – number of results per query

**get\_query\_results** (*query\_id*)

Returns the corresponding RetrievalResults object for a given query.

**Parameters** **query\_id** – queryID

**Return type** `nordlys.core.retrieval.retrieval_results.RetrievalResults`

**get\_results** ()

Returns all results.

**Returns** a dict with queryIDs as keys and RetrievalResults object as values

**load\_file** (*file\_name*, *remap\_by\_exp=False*)

Loads a TREC runfile.

**Parameters**

- **file\_name** – name of the run file
- **remap\_by\_exp** – whether scores are to be converted from the log-domain by taking their exp (default: False)

**normalize** ()

Normalizes the retrieval scores such that they sum up to one for each query.

**print\_stat** ()

Prints simple statistics.

`nordlys.core.eval.trec_run.arg_parser` ()

`nordlys.core.eval.trec_run.main` (*args*)

## nordlys.core.ml package

### Machine learning

The machine learning package is connected to [Scikit-learn](#) and can be used for learning-to-rank and classification purposes.

### Usage

For information on how to use the package from command line and set the configuration, read ML usage

### Data format

The file format of the training and test files is *json*. Each instance is presented as a dictionary, consist of the following elements:

- **ID**: Instance id
- **Target**: The target value of the instance.
- **Features**: All the features, presented in key-value format. Note that all the instances should have the same set of features.

- **Properties:** All meta-data about the instance; e.g. query ID, or content.

Below is an excerpt from a json data file:

```
{
  "0": {
    "properties": {
      "query": "papaqui soccer",
      "entity": "<dbpedia:Soccer_(1985_video_game)>"
    },
    "target": "0",
    "features": {
      "feat1": 1,
      "feat2": 0,
      "feat3": 25
    }
  },
  "1": {}
}
```

---

**Note:** The sample files for using the ML package are provided under `data/ml_sample/` folder.

---

**Note:** Currently we provide support for Random Forest(RF) and Gradient Boosted Regression Trees (GBRT).

---

## Submodules

### nordlys.core.ml.cross\_validation module

#### Cross Validation

Cross-validation support.

We assume that instances (i) are uniquely identified by an instance ID and (ii) they have id and score properties. We access them using the Instances class.

**Authors** Faegheh Hasibi, Krisztian Balog

**class** nordlys.core.ml.cross\_validation.**CrossValidation** (*k*, *instances*, *callback\_train*, *callback\_test*)

Bases: `object`

**Class attributes:**

**fold:** dict of folds (1..k) with a dict {"training": [list of instance\_ids]}, {"testing": [list of instance\_ids]}

**Parameters**

- **k** – number of folds
- **instances** – Instances object
- **callback\_train** – Callback function for training model
- **callback\_test** – Callback function for applying model

**create\_folds** (*group\_by=None*)

Creates folds for the data set.

**Parameters** **group\_by** – property to group by (instance\_id by default)

**get\_folds** (*filename=None, group\_by=None*)

Loads folds from file or generates them if the file doesn't exist.

**Parameters**

- **filename** –
- **k** – number of folds

**Returns**

**get\_instances** (*i, mode, property=None*)

Returns instances from the given fold *i* in [0..k-1].

**Parameters**

- **i** – fold number
- **mode** – training or testing

:return Instances object

**load\_folds** (*filename*)

Loads previously created folds from (JSON) file.

**run** ()

Runs cross-validation.

**save\_folds** (*filename*)

Saves folds to (JSON) file.

`nordlys.core.ml.cross_validation.main()`

## nordlys.core.ml.instance module

### Instance

Instance class.

#### Features:

- This class supports different features for an instance.
- The features, together with the id of the instance, will be used in machine learning algorithms.
- All Features are stored in a dictionary, where keys are feature names (self.features).

#### Instance properties:

- Properties are additional side information of an instance (e.g. query\_id, entity\_id, ...).
- properties are stored in a dictionary (self.properties).

This is the base instance class. Specific type of instances can inherit from class and add more properties to the base class.

**Author** Faegheh Hasibi

**class** `nordlys.core.ml.instance.Instance` (*id, features=None, target='0', properties=None*)

Bases: `object`

**Class attributes:** `ins_id`: (string) features: a dictionary of feature names and values `target`: (string) target id or class properties: a dictionary of property names and values

**add\_feature** (*feature*, *value*)

Adds a new feature to the features.

**Parameters** `feature` – (string), feature name

:param value

**add\_property** (*property*, *value*)

Adds a new property to the properties.

**Parameters** `property` – (string), property name

:param value

**features**

**classmethod** `from_json` (*ins\_id*, *fields*)

Reads an instance in JSON format and generates Instance.

**Parameters**

- `ins_id` – instance id
- `fields` – A dictionary of fields

:return (ml.Instance)

**get\_feature** (*feature*)

Returns the value of a given feature.

:param feature :return value

**get\_property** (*property*)

Returns the value of a given property.

:param property :return value

**id**

**properties**

**to\_json** (*file\_name=None*)

Converts instance to the JSON format.

**Parameters** `file_name` – (string)

:return JSON dump of the instance.

**to\_libsvm** (*features*, *qid\_prop=None*)

Converts instance to the Libsvm format. - RankLib format:

<target> qid:<qid> <feature>:<value> ... # <info>

- Example: 3 qid:1 1:1 2:1 3:0 4:0.2 5:0 # 1A

NOTE: the property used for qid(`qid_prop`) should hold integers

**Parameters**

- `features` – the list of features that should be in the output
- `qid_prop` – property to be used as qid

:return str, instance in the rankLib format.

**to\_str** (*feature\_set=None*)  
Converts instances to string.

**Parameters** **feature\_set** – features to be included in the output format

**Return (string)** **tab separated string** ins\_id target ftr\_1 ftr\_2 ... ftr\_n properties

`nordlys.core.ml.instance.main()`

## nordlys.core.ml.instances module

### Instances

Instances used for Machine learning algorithms.

- Manages a set of Instance objects
- **Loads instance-data from JSON or TSV files**
  - When using TSV, instance properties, target, and features are loaded from separate files
- Generates a list of instances in JSON or RankLib format

**Authors** Faegheh Hasibi, Krisztian Balog

**class** `nordlys.core.ml.instances.Instances` (*instances=None*)

Bases: `object`

**Class attributes:** instances: Instance objects stored in a dictionary indexed by instance id

**Parameters** **instances** – instances in a list or dict - if list then list index is used as the instance ID - if dict then the key is used as the instance ID

**add\_features\_from\_tsv** (*tsv\_file, features*)

**add\_instance** (*instance*)

Adds an Instance object to the list of instances.

**Parameters** **instance** – Instance object

**add\_properties\_from\_tsv** (*tsv\_file, properties*)

**add\_qids** (*prop*)

Generates (integer) q\_id-s (for libsvm) based on a given (non-integer) property. It assigns a unique integer value to each different value for that property.

**Parameters** **prop** – name of the property.

**Returns**

**add\_target\_from\_tsv** (*tsv\_file*)

**append\_instances** (*ins\_list*)

Appends the list of Instances objects.

**Parameters** **ins\_list** – list of Instance objects

**classmethod** **from\_json** (*json\_file*)

Loads instances from a JSON file.

**Parameters** **json\_file** – (string)

:return Instances object

**get\_all** ()

Returns list of all instances.

**get\_all\_ids** ()

Returns list of all instance ids.

**get\_instance** (*instance\_id*)

Returns an instance by instance id.

**Parameters** *instance\_id* – (string)

**Returns** Instance object

**group\_by\_property** (*property*)

Groups instances by a given property.

:param *property* :return a dictionary of instance ids {id:[ml.Instance, ...], ... }

**to\_json** (*json\_file=None*)

Converts all instances to JSON and writes it to the file

**Parameters** *json\_file* – (string)

**Returns** JSON dump of all instances.

**to\_libsvm** (*file\_name=None, qid\_prop=None*)

Converts all instances to the LibSVM format and writes them to the file. - Libsvm format:

<line> .=. <target> qid:<qid> <feature>:<value> ... # <info> <target> .=. <float> <qid> .=.  
<positive integer> <feature> .=. <positive integer> <value> .=. <float> <info> .=. <string>

- Example: 3 qid:1 1:1 2:1 3:0 4:0.2 5:0 # 1A

#### NOTES:

- The property used for qid(*qid\_prop*) should hold integers
- For pointwise algorithms, we use instance id for qid
- Lines in the RankLib input have to be sorted by increasing qid.

#### Parameters

- **file\_name** – File to write libsvm format of instances.
- **qid\_prop** – property to be used as qid. If none,

**to\_str** (*file\_name=None*)

Converts instances to string and write them to the given file. :param *file\_name* :return: String format of instances

**to\_treceval** (*file\_name, qid\_prop='qid', docid\_prop='en\_id'*)

Generates a TREC style run file - If there is an entity ranked more than once for the same query, the one with higher score is kept.

#### Parameters

- **file\_name** – File to write TREC file
- **qid\_prop** – Name of instance property to be used as query ID (1st column)
- **docid\_prop** – Name of instance property to be used as document ID (3rd column)

nordlys.core.ml.instances.main (*args*)

## nordlys.core.ml.ml module

### Machine leaning

The command-line application for general-purpose machine learning.

### Usage

```
python -m nordlys.core.ml.ml <config_file>
```

### Config parameters

- **training\_set**: nordlys ML instance file format (MIFF)
- **test\_set**: nordlys ML instance file format (MIFF); if provided then it's always used for testing. Can be left empty if cross-validation is used, in which case the remaining split is used for testing.
- **cross\_validation**:
  - **k**: number of folds (default: 10); use -1 for leave-one-out
  - **split\_strategy**: name of a property (normally query-id for IR problems). If set, the entities with the same value for that property are kept in the same split. if not set, entities are randomly distributed among splits.
  - **splits\_file**: JSON file with splits (instance\_ids); if the file is provided it is used, otherwise it's generated
  - **create\_splits**: if True, creates the CV splits. Otherwise loads the splits from "split\_file" parameter.
- **model**: ML model, currently supported values: rf, gbrt
- **category**: [regression | classification], default: "regression"
- **parameters**: **dict with parameters of the given ML model**
  - **If GBRT**:
    - \* **alpha**: learning rate, default: 0.1
    - \* **tree**: number of trees, default: 1000
    - \* **depth**: max depth of trees, default: 10% of number of features
  - **If RF**:
    - \* **tree**: number of trees, default: 1000
    - \* **maxfeat**: max features of trees, default: 10% of number of features
- **model\_file**: the model is saved to this file
- **load\_model**: if True, loads the model
- **feature\_imp\_file**: Feature importance is saved to this file
- **output\_file**: where output is written; default output format: TSV with with instance\_id and (estimated) target

## Example config

```
{
  "model": "gbrt",
  "category": "regression",
  "parameters":{
    "alpha": 0.1,
    "tree": 10,
    "depth": 5
  },
  "training_set": "path/to/train.json",
  "test_set": "path/to/test.json",
  "model_file": "path/to/model.txt",
  "output_file": "path/to/output.json",
  "cross_validation":{
    "create_splits": true,
    "splits_file": "path/to/splits.json",
    "k": 5,
    "split_strategy": "q_id"
  }
}
```

---

**Authors** Faegheh Hasibi, Krisztian Balog

**class** nordlys.core.ml.ml.**ML** (*config*)

Bases: object

**analyse\_features** (*model, feature\_names*)

Ranks features based on their importance. Scikit uses Gini score to get feature importances.

**Parameters**

- **model** – trained model
- **feature\_names** – list of feature names

**apply\_model** (*instances, model*)

Applies model on a given set of instances.

**Parameters**

- **instances** – Instances object
- **model** – trained model

**Returns** Instances

**gen\_model** (*num\_features=None*)

Reads parameters and generates a model to be trained.

**Parameters** **num\_features** – int, number of features

:return untrained ranker/classifier

**output** (*instances*)

Writes results to output file.

**Parameters** **instances** – Instances object

**run** ()

**train\_model** (*instances*)

Trains model on a given set of instances.

**Parameters** *instances* – Instances object

**Returns** the learned model

```
nordlys.core.ml.ml.arg_parser()
```

```
nordlys.core.ml.ml.main(args)
```

## nordlys.core.retrieval package

### Retrieval

The retrieval package provides basic indexing and scoring functionality based on Elasticsearch (v2.3). It can be used both for documents and for entities (as the latter are represented as fielded documents).

### Indexing

Indexing can be done by directly reading the content of documents. The *toy\_indexer* module provides a toy example.

When the content of documents is stored in MongoDB (e.g., for DBpedia entities), use the *indexer\_mongo* module for indexing. For further details on how this module can be used, see *indexer\_dbpedia*.

For indexing Dbpedia entities, we read the content of entiteis form MongoDB aFor DBpedia entities, we store them on MongoDB and .. todo:: Explain indexing (representing entities as fielded documents, mongo to elasticsearch)

### Notes

- To speed up indexing, use *add\_docs\_bulk()*. The optimal number of documents to send in a single bulk depends on the size of documents; you need to figure it out experimentally.
- We strongly recommend using the default Elasticsearch similarity (currently BM25) for indexing. (Other similarity functions may be also used; in that case the similarity function can updated after indexing.)
- Our default setting is *not* to store term positions in the index (for efficiency considerations).

### Retrieval

Retrieval is done in two stages:

- *First pass*: The top N documents are retrieved using Elastic’s default search method
- *Second pass*: The (expensive) scoring of the top N documents is performed (implemented in the Nordlys)

Nordlys currently supports the following models for second pass retrieval:

- Language modelling (LM) [1]
- Mixture of Language Modesl (MLM) [2]
- Probabilistic Model for Semistructured Data (PRMS) [3]

Check out *scorer* module to get inspiration for implementing a new retrieval model.

### Command line usage

See `nordlys.core.retrieval.retrieval`

### Notes

- Always use a `ElasticCache` object (instead of `Elastic`) for getting stats from the index. This class stores index stats in the memory, which highly benefits efficiency.
- We recommend to create a new `ElasticCache` object for each query. This way, you will make effiecnt of your machine's memory.

---

[1] Jay M Ponte and W Bruce Croft. 1998. *A Language modeling approach to information retrieval*. In Proc. of SIGIR '98.

[2] Paul Ogilvie and Jamie Callan. 2003. *Combining document representations for known-item search*. Proc. of SIGIR '03.

[3] Jinyoung Kim, Xiaobing Xue, and W Bruce Croft. 2009. *A probabilistic retrieval model for semistructured data*. In Proc. of ECIR '09.

### Submodules

#### `nordlys.core.retrieval.elastic` module

#### Elastic

Utility class for working with Elasticsearch. This class is to be instantiated for each index.

#### Indexing usage

To create an index, first you need to define field mappings and then build the index. The sample code for creating an index is provided at `nordlys.core.retrieval.toy_indexer`.

#### Retrieval usage

The following statistics can be obtained from this class:

- Number of documents: `Elastic.num_docs()`
- Number of fields: `Elastic.num_fields()`
- Document count: `Elastic.doc_count()`
- Collection length: `Elastic.coll_length()`
- Average length: `Elastic.avg_len()`
- Document length: `Elastic.doc_length()`
- Document frequency: `Elastic.doc_freq()`
- Collection frequency: `Elastic.coll_term_freq()`

- Term frequencies: `Elastic.term_freqs()`

## Efficiency considerations

- For efficiency reasons, we do not store term positions during indexing. To store them, see the corresponding mapping functions `Elastic.analyzed_field()`, `Elastic.notanalyzed_searchable_field()`.
- Use `ElasticCache` for getting index statistics. This module caches the statistics into memory and boosts efficiency.
- Mind that `ElasticCache` does not empty the cache!

**Authors** Faegheh Hasibi, Krisztian Balog

**class** nordlys.core.retrieval.elastic.**Elastic** (*index\_name*)

Bases: `object`

**ANALYZER\_STOP** = 'stop\_en'

**ANALYZER\_STOP\_STEM** = 'english'

**BM25** = 'BM25'

**DOC\_TYPE** = 'doc'

**FIELD\_CATCHALL** = 'catchall'

**FIELD\_ELASTIC\_CATCHALL** = '\_all'

**SIMILARITY** = 'sim'

**add\_doc** (*doc\_id*, *contents*)

Adds a document with the specified contents to the index.

### Parameters

- **doc\_id** – document ID
- **contents** – content of document

**add\_docs\_bulk** (*docs*)

Adds a set of documents to the index in a bulk.

**Parameters docs** – dictionary {doc\_id: doc}

**analyze\_query** (*query*, *analyzer*='stop\_en')

Analyzes the query.

### Parameters

- **query** – raw query
- **analyzer** – name of analyzer

**static analyzed\_field** (*analyzer*='stop\_en')

Returns the mapping for analyzed fields.

For efficiency considerations, term positions are not stored. To store term positions, change "term\_vector": "with\_positions\_offsets"

**Parameters analyzer** – name of the analyzer; valid options: [ANALYZER\_STOP, ANALYZER\_STOP\_STEM]

**avg\_len** (*field*)

Returns average length of a field in the collection.

**coll\_length** (*field*)

Returns length of field in the collection.

**coll\_term\_freq** (*term, field, tv=None*)

Returns collection term frequency for the given field.

**create\_index** (*mappings, model='BM25', model\_params=None, force=False*)

Creates index (if it doesn't exist).

**Parameters**

- **mappings** – field mappings
- **model** – name of elastic search similarity
- **model\_params** – name of elastic search similarity
- **force** – forces index creation (overwrites if already exists)

**delete\_index** ()

Deletes an index.

**doc\_count** (*field*)

Returns number of documents with at least one term for the given field.

**doc\_freq** (*term, field, tv=None*)

Returns document frequency for the given term and field.

**doc\_length** (*doc\_id, field*)

Returns length of a field in a document.

**get\_doc** (*doc\_id, fields=None, source=True*)

Gets a document from the index based on its ID.

**Parameters**

- **doc\_id** – document ID
- **fields** – list of fields to return (default: all)
- **source** – return document source as well (default: yes)

**get\_field\_stats** (*field*)

Returns stats of the given field.

**get\_fields** ()

Returns name of fields in the index.

**get\_mapping** ()

Returns mapping definition for the index.

**get\_settings** ()

Returns index settings.

**static notanalyzed\_field** ()

Returns the mapping for not-analyzed fields.

**static notanalyzed\_searchable\_field** ()

Returns the mapping for not-analyzed fields.

**num\_docs** ()

Returns the number of documents in the index.

**num\_fields** ()

Returns number of fields in the index.

**search** (*query*, *field*, *num=100*, *fields\_return=""*, *start=0*)

Searches in a given field using the similarity method configured in the index for that field.

#### Parameters

- **query** – query string
- **field** – field to search in
- **num** – number of hits to return (default: 100)
- **fields\_return** – additional document fields to be returned
- **start** – starting offset (default: 0)

**Returns** dictionary of document IDs with scores

**search\_complex** (*body*, *num=10*, *fields\_return=""*, *start=0*)

Supports complex structured queries, which are sent as a `body` field in Elastic search. For detailed information on formulating structured queries, see the [official instructions](#). Below is an example to search in two particular fields that each must contain a specific term.

#### Example

```
# [explanation of the query]
term_1 = "hello"
term_2 = "world"
body = {
  "query": {
    "bool": {
      "must": [
        {
          "match": {"title": term_1}
        },
        {
          "match_phrase": {"content": term_2}
        }
      ]
    }
  }
}
```

#### Parameters

- **body** – query body
- **field** – field to search in
- **num** – number of hits to return (default: 100)
- **fields\_return** – additional document fields to be returned
- **start** – starting offset (default: 0)

**Returns** dictionary of document IDs with scores

**term\_freq** (*doc\_id*, *field*, *term*)

Returns frequency of a term in a given document and field.

**term\_freqs** (*doc\_id*, *field*, *tv=None*)

Returns term frequencies of all terms for a given document and field.

**update\_similarity** (*model='BM25', params=None*)

Updates the similarity function “sim”, which is fixed for all index fields.

The method and param should match elastic settings: <https://www.elastic.co/guide/en/elasticsearch/reference/2.3/index-modules-similarity.html>

#### Parameters

- **model** – name of the elastic model
- **params** – dictionary of params based on elastic

## nordlys.core.retrieval.elastic\_cache module

### Elastic Cache

This is a cache for elastic index stats; a layer between an index and retrieval. The statistics (such as document and term frequencies) are first read from the index and stay in the memory for further usages.

### Usage hints

- Only one instance of Elastic cache needs to be created.
- If running out of memory, you need to create a new object of ElasticCache.
- The class also caches termvectors. To further boost efficiency, you can load term vectors for multiple documents using `ElasticCache.multi_termvector()`.

**Author** Faegheh Hasibi

**class** nordlys.core.retrieval.elastic\_cache.**ElasticCache** (*index\_name*)

Bases: `nordlys.core.retrieval.elastic.Elastic`

**avg\_len** (*field*)

Returns average length of a field in the collection.

**coll\_length** (*field*)

Returns length of field in the collection.

**coll\_term\_freq** (*term, field, tv=None*)

Returns collection term frequency for the given field.

**doc\_count** (*field*)

Returns number of documents with at least one term for the given field.

**doc\_freq** (*term, field, tv=None*)

Returns document frequency for the given term and field.

**doc\_length** (*doc\_id, field*)

Returns length of a field in a document.

**multi\_termvector** (*doc\_ids, field, batch=50*)

Returns term vectors for a given document and field.

**num\_docs** ()

Returns the number of documents in the index.

**num\_fields** ()

Returns number of fields in the index.

**term\_freq** (*doc\_id, field, term*)

Returns frequency of a term in a given document and field.

**term\_freqs** (*doc\_id, field, tv=None*)

Returns term frequencies for a given document and field.

## nordlys.core.retrieval.indexer\_mongo module

### Mongo Indexer

This class is a tool for creating an index from a Mongo collection.

To use this class, you need to implement `callback_get_doc_content()` function. See `indexer_fsdm` for an example usage of this class.

**Author** Faegheh Hasibi

```
class nordlys.core.retrieval.indexer_mongo.IndexerMongo (index_name,           map-
                                                    pings,           collection,
                                                    model='BM25')
```

Bases: `object`

**build** (*callback\_get\_doc\_content, bulk\_size=1000*)

Builds the DBpedia index from the mongo collection.

To speedup indexing, we index documents as a bulk. There is an optimum value for the bulk size; try to figure it out.

#### Parameters

- **callback\_get\_doc\_content** – a function that get a documet from mongo and return the content for indexing
- **bulk\_size** – Number of documents to be added to the index as a bulk

## nordlys.core.retrieval.retrieval module

### Retrieval

Console application for general-purpose retrieval.

### Usage

```
python -m nordlys.services.er -c <config_file> -q <query>
```

If `-q <query>` is passed, it returns the results for the specified query and prints them in terminal.

### Config parameters

- **index\_name**: name of the index,
- **first\_pass**:
  - **1st\_num\_docs**: number of documents in first-pass scoring (default: 100)

- **field**: field used in first pass retrieval (default: Elastic.FIELD\_CATCHALL)
- **fields\_return**: comma-separated list of fields to return for each hit (default: “”)
- **num\_docs**: number of documents to return (default: 100)
- **start**: starting offset for ranked documents (default:0)
- **model**: name of retrieval model; accepted values: [lm, mlm, prms] (default: lm)
- **field**: field name for LM (default: catchall)
- **fields: single field name for LM (default: catchall)** list of fields for PRMS (default: [catchall]) dictionary with fields and corresponding weights for MLM (default: {catchall: 1})
- **smoothing\_method**: accepted values: [jm, dirichlet] (default: dirichlet)
- **smoothing\_param**: value of lambda or mu; accepted values: [float or “avg\_len”], (jm default: 0.1, dirichlet default: 2000)
- **query\_file**: name of query file (JSON),
- **output\_file**: name of output file,
- **run\_id**: run id for TREC output

### Example config

```
{
  "index_name": "dbpedia_2015_10",
  "first_pass": {
    "1st_num_docs": 1000
  },
  "model": "prms",
  "num_docs": 1000,
  "smoothing_method": "dirichlet",
  "smoothing_param": 2000,
  "fields": ["names", "categories", "attributes", "similar_entity_names", "related_
  ↪entity_names"],
  "query_file": "path/to/queries.json",
  "output_file": "path/to/output.txt",
  "run_id": "test"
}
```

**Authors** Krisztian Balog, Faegheh Hasibi

**class** nordlys.core.retrieval.retrieval.**Retrieval** (*config*)

Bases: object

**FIELDDED\_MODELS** = set(['mlm', 'prms'])

**LM\_MODELS** = set(['lm', 'mlm', 'prms'])

**batch\_retrieval** ()

Scores queries in a batch and outputs results.

**static check\_config** (*config*)

Checks config parameters and sets default values.

**retrieve** (*query*, *scorer=None*)

Scores documents for the given query.

```
trec_format (results, query_id, max_rank=100)
    Outputs results in TREC format
```

```
nordlys.core.retrieval.retrieval.arg_parser ()
```

```
nordlys.core.retrieval.retrieval.get_config ()
```

```
nordlys.core.retrieval.retrieval.main (args)
```

## nordlys.core.retrieval.retrieval\_results module

### Retrieval Results

Result list representation.

- for each hit it holds score and both internal and external doc\_ids

**Authors** Faegheh Hasibi, Krisztian Balog

```
class nordlys.core.retrieval.retrieval_results.RetrievalResults (scores={},  
                                                             query=None)
```

Bases: `object`

Class for storing retrieval scores for a given query.

```
append (doc_id, score)
    Adds document to the result list
```

```
classmethod elastic_to_retrieval (res, query=None)
    Converts elastic search results to retrieval results.
```

```
get_score (doc_id)
    Returns the score of a document (or None if it's not in the list).
```

```
get_scores_sorted ()
    Returns all results sorted by score
```

```
num_docs ()
    Returns the number of documents in the result list.
```

**query**

```
write_trec_format (query_id, run_id, out, max_rank=100)
    Outputs results in TREC format
```

## nordlys.core.retrieval.scorer module

### Scorer

Various retrieval models for scoring a individual document for a given query.

**Authors** Faegheh Hasibi, Krisztian Balog

```
class nordlys.core.retrieval.scorer.Scorer (elastic, query, params)
    Bases: object
```

Base scorer class.

```
SCORER_DEBUG = 0
```

**static** `get_scorer` (*elastic, query, config*)  
 Returns Scorer object (Scorer factory).

**Parameters**

- **elastic** – Elastic object
- **query** – raw query (to be analyzed)
- **config** – dict with models parameters

**class** `nordlys.core.retrieval.scorer.ScorerLM` (*elastic, query, params*)  
 Bases: `nordlys.core.retrieval.scorer.Scorer`

Language Model (LM) scorer.

**DIRICHLET** = 'dirichlet'

**JM** = 'jm'

**static** `get_dirichlet_prob` (*tf\_t\_d, len\_d, tf\_t\_C, len\_C, mu*)  
 Computes Dirichlet-smoothed probability.  $P(t|\theta_d) = [tf(t, d) + \mu P(t|C)] / [|d| + \mu]$

**Parameters**

- **tf\_t\_d** –  $tf(t, d)$
- **len\_d** –  $|d|$
- **tf\_t\_C** –  $tf(t, C)$
- **len\_C** –  $|C| = \sum_{d \in C} |d|$
- **mu** –  $\mu$

**Returns** Dirichlet-smoothed probability

**static** `get_jm_prob` (*tf\_t\_d, len\_d, tf\_t\_C, len\_C, lambda*)  
 Computes JM-smoothed probability.  $p(t|\theta_d) = [(1-\lambda) tf(t, d)/|d|] + [\lambda tf(t, C)/|C|]$

**Parameters**

- **tf\_t\_d** –  $tf(t, d)$
- **len\_d** –  $|d|$
- **tf\_t\_C** –  $tf(t, C)$
- **len\_C** –  $|C| = \sum_{d \in C} |d|$
- **lambda** –  $\lambda$

**Returns** JM-smoothed probability

**get\_lm\_term\_prob** (*doc\_id, field, t, tf\_t\_d\_f=None, tf\_t\_C\_f=None*)  
 Returns term probability for a document and field.

**Parameters**

- **doc\_id** – document ID
- **field** – field name
- **t** – term

**Returns**  $P(t|d_f)$

**get\_lm\_term\_probs** (*doc\_id, field*)  
 Returns probability of all query terms for a document and field; i.e.  $p(t|\theta_d)$

**Parameters**

- **doc\_id** – document ID
- **field** – field name

**Returns** dictionary of terms with their probabilities

**score\_doc** (*doc\_id*)

Scores the given document using LM.  $p(q|t\theta_d) = \sum \log(p(t|\theta_d))$

**Parameters** **doc\_id** – document id

**Returns** LM score

**class** nordlys.core.retrieval.scorer.**ScorerMLM** (*elastic, query, params*)

Bases: *nordlys.core.retrieval.scorer.ScorerLM*

Mixture of Language Model (MLM) scorer.

**Implemented based on:** Ogilvie, Callan. Combining document representations for known-item search. SIGIR 2003.

**get\_mlm\_term\_prob** (*doc\_id, t*)

Returns MLM probability for the given term and field-weights.  $p(t|\theta_d) = \sum (\mu_f * p(t|\theta_{d_f}))$

**Parameters**

- **lucene\_doc\_id** – internal Lucene document ID
- **t** – term

**Returns**  $P(t|\theta_d)$

**get\_mlm\_term\_probs** (*doc\_id*)

Returns probability of all query terms for a document; i.e.  $p(t|\theta_d)$

**Parameters** **doc\_id** – internal Lucene document ID

**Returns** dictionary of terms with their probabilities

**score\_doc** (*doc\_id*)

Scores the given document using MLM model.  $p(q|t\theta_d) = \sum \log(p(t|\theta_d))$

**Parameters** **doc\_id** – document ID

**Returns** MLM score of document and query

**class** nordlys.core.retrieval.scorer.**ScorerPRMS** (*elastic, query, params*)

Bases: *nordlys.core.retrieval.scorer.ScorerLM*

PRMS scorer.

**get\_mapping\_prob** (*t, coll\_termfreq\_fields=None*)

**Computes PRMS field mapping probability.**  $p(ft) = P(tf)P(f) / \sum_{f'} (P(t|C_{\{f'_c\}})P(f'))$

**Parameters**

- **t** – str
- **coll\_termfreq\_fields** – {field: freq, ... }

**Returns** a dictionary {field: prms\_prob, ... }

**get\_mapping\_probs** ()

Gets (cached) mapping probabilities for all query terms.

`get_total_field_freq()`

Returns total occurrences of all fields

`score_doc(doc_id)`

Scores the given document using PRMS model.

### Parameters

- `doc_id` – document id
- `lucene_doc_id` – internal Lucene document ID

**Returns** float, PRMS score of document and query

## nordlys.core.retrieval.toy\_indexer module

### Toy Indexer

Toy indexing example for testing purposes.

**Authors** Krisztian Balog, Faegheh Hasibi

`nordlys.core.retrieval.toy_indexer.main()`

## nordlys.core.storage package

This is the (intro of the) storage package doc. It's in the `__init__.py` of the package, as an example for the other packages, and IT SHOULD BE cleaned up without all of this ;).

No idea what text it could contain :)

It might contain, e.g., items in this way

- This
- That

and also code snippets like usual, using as usual these double colons and the indented code:

```
mongo = Mongo(host, db, collection)
```

And that's it.

(In the following I included all the up-to-now available modules, sectioned by no criteria in particular.)

### Subpackages

#### nordlys.core.storage.parser package

### Submodules

#### nordlys.core.storage.parser.nt\_parser module

### NTriples Parser

NTriples parser with URI prefixing

**Author** Krisztian Balog

**class** nordlys.core.storage.parser.nt\_parser.NTParser

Bases: `object`

NTriples parser class

**parse\_file** (*filename, triplehandler*)

Parses file and calls callback function with the parsed triple

**class** nordlys.core.storage.parser.nt\_parser.Triple (*prefix=None*)

Bases: `object`

Representation of a Triple to be used by the rdflib NTriplesParser.

**object** ()

**object\_prefixed** ()

**predicate** ()

**predicate\_prefixed** ()

**subject** ()

**subject\_prefixed** ()

**triple** (*s, p, o*)

Assign current triple object

#### Parameters

- **s** – subject
- **p** – predicate
- **o** – object

**class** nordlys.core.storage.parser.nt\_parser.TripleHandler

Bases: `object`

This is an abstract class

**triple\_parsed** (*triple*)

This method is called each time a triple is parsed, with the triple as parameter.

**class** nordlys.core.storage.parser.nt\_parser.TripleHandlerPrinter

Bases: `nordlys.core.storage.parser.nt_parser.TripleHandler`

Example triple handler that only prints whatever it received.

**triple\_parsed** (*triple*)

This method is called each time a triple is parsed, with the triple as parameter.

nordlys.core.storage.parser.nt\_parser.**main** (*argv*)

## nordlys.core.storage.parser.uri\_prefix module

### URI Prefixing

URI prefixing.

**Author** Krisztian Balog

**class** nordlys.core.storage.parser.uri\_prefix.**URIPrefix** (*prefix\_file*='/home/docs/checkouts/readthedocs.org/u  
Bases: `object`

**get\_prefixed** (*uri*, *angle\_brackets*=*True*)

nordlys.core.storage.parser.uri\_prefix.**convert\_txt\_to\_json** (*txt\_file*,  
*json\_file*='/home/docs/checkouts/readthedocs.o

Convert prefixes txt file to json.

This has to be done only once. And only in case there is no .json file, or any changes done in .txt.

## Submodules

### nordlys.core.storage.mongo module

#### Mongo

Tools for working with MongoDB.

**Authors** Krisztian Balog, Faegheh Hasibi

**class** nordlys.core.storage.mongo.**Mongo** (*host*, *db*, *collection*)

Bases: `object`

Manages the MongoDB connection and operations.

**ID\_FIELD** = `'_id'`

**add** (*doc\_id*, *contents*)

Adds a document or replaces the contents of an entire document.

**append\_dict** (*doc\_id*, *field*, *dictkey*, *value*)

Appends the value to a given field that stores a dict. If the dictkey is already in use, the value stored there will be overwritten.

#### Parameters

- **doc\_id** – document id
- **field** – field
- **dictkey** – key in the dictionary
- **value** – value to be increased by

**append\_list** (*doc\_id*, *field*, *value*)

Appends the value to a given field that stores a list. If the field does not exist yet, it will be created. The value should be a list.

#### Parameters

- **doc\_id** – document id
- **field** – field
- **value** – list, a value to be appended to the current list

**append\_set** (*doc\_id*, *field*, *value*)

Adds a list of values to a set. If the field does not exist yet, it will be created. The value should be a list.

#### Parameters

- **doc\_id** – document id

- **field** – field
- **value** – list, a value to be appended to the current list

**drop** ()

Deletes the contents of the given collection (including indices).

**find\_all** (*no\_timeout=False*)

Returns a Cursor instance that allows us to iterate over all documents.

**find\_by\_id** (*doc\_id*)

Returns unescaped document content for a given document id.

**get\_num\_docs** ()

Returns total number of documents in the mongo collection.

**inc** (*doc\_id, field, value*)

Increments the value of a specified field.

**inc\_in\_dict** (*doc\_id, field, dictkey, value=1*)

Increments a value that is inside a dict.

#### Parameters

- **doc\_id** – document id
- **field** – field
- **dictkey** – key in the dictionary
- **value** – value to be increased by

**static print\_doc** (*doc*)

**set** (*doc\_id, field, value*)

Sets the value of a given document field (overwrites previously stored content).

**static unescape** (*s*)

Unescapes string.

**static unescape\_doc** (*mdoc*)

Unescapes document content.

`nordlys.core.storage.mongo.main()`

## nordlys.core.storage.nt2mongo module

## nordlys.core.utils package

@author: Krisztian Balog

## Submodules

### nordlys.core.utils.entity\_utils module

### nordlys.core.utils.file\_utils module

## File Utils

Utility methods for file handling.

**Authors** Krisztian Balog, Faegheh Hasibi

**class** nordlys.core.utils.file\_utils.**FileUtils**

Bases: `object`

**static dump\_tsv** (*file\_name*, *data*, *header=None*, *append=False*)

Dumps the data in tsv format.

### Parameters

- **file\_name** – name of file
- **data** – list of list
- **header** – list of headers
- **append** – if True, appends the data to the existing file

**static load\_config** (*config*)

Loads config file/dictionary.

**Parameters** **config** – json file or a dictionary

**Returns** config dictionary

**static open\_file\_by\_type** (*file\_name*, *mode='r'*)

Opens file (gz/text) and returns the handler.

**Parameters** **file\_name** – NTriples file

**Returns** handler to the file

**static read\_file\_as\_list** (*filename*)

Reads in non-empty lines from a textfile (which may be gzipped/bz2ed) and returns it as a list.

**Parameters** **filename** –

nordlys.core.utils.file\_utils.**main**()

## nordlys.core.utils.logging\_utils module

### Logging Utils

Utility methods for logging.

**Author** Heng Ding

**class** nordlys.core.utils.logging\_utils.**PrintHandler** (*logging\_level*)

Bases: `object`

Handler for elastic prints

**class** nordlys.core.utils.logging\_utils.**RequestHandler** (*logging\_path*)

Bases: `object`

Handler for elastic request

## nordlys.logic package

### Services

All modules in this package serving for services layer.

### Subpackages

#### nordlys.logic.el package

### Entity linking

This package is the implementation of entity linking.

### Submodules

#### nordlys.logic.el.cmns module

### Commonness Entity Linking Approach

Class for commonness entity linking approach

**Author** Faegheh Hasibi

```
class nordlys.logic.el.cmns.Cmns (query, entity, threshold=None, cmns_th=0.1)
```

Bases: `object`

```
disambiguate ()
```

Selects only one entity per mention.

```
Return [{"mention": xx, "entity": yy, "score": zz}, ...] #dictionary {mention: (en_id, score), ..}
```

```
link ()
```

Links the query to the entity.

```
dictionary {mention: (en_id, score), ..}
```

```
rank_ens ()
```

Detects mention and rank entities for each mention

```
nordlys.logic.el.cmns.main (args)
```

#### nordlys.logic.el.el\_utils module

### EL Utils

Utility methods for entity linking.

Author: Faegheh Hasibi

```
nordlys.logic.el.el_utils.is_name_entity (en_id)
```

Returns true if the entity is considered as proper name entity.

`nordlys.logic.el.el_utils.load_kb_snapshot` (*kb\_file*)

Loads DBpedia Snapshot of proper name entities (used for entity linking).

`nordlys.logic.el.el_utils.to_elq_eval` (*annotations, output\_file*)

Write entity annotations to ELQ evaluation format.

**Parameters** `linked_ens` – {qid:[{"mention":xx, "entity": yy, "score":zz}, ..], ..}

## nordlys.logic.el.greedy module

Generative model for interpretation set finding

@author: Faegheh Hasibi

**class** `nordlys.logic.el.greedy.Greedy` (*score\_th*)

Bases: `object`

**create\_interpretations** (*query\_inss*)

Groups CER instances as interpretation sets.

**Return list of interpretations, where each interpretation is a dictionary** {`mention` (*en\_id*, *score*), ..}

**disambiguate** (*inss*)

Takes instances and generates set of entity linking interpretations.

**Parameters** `inss` – Instances object

**Returns** sets of interpretations [{`mention`: (*en\_id*, *score*), ..}, ...]

**is\_overlapping** (*mentions*)

Checks whether the strings of a set overlapping or not. i.e. if there exists a term that appears twice in the whole set.

**E.g.** {"the", "music man"} is not overlapping {"the", "the man", "music"} is overlapping.

NOTE: If a query is "yxxz" the mentions {"yx", "xz"} and {"yx", "x"} are overlapping.

**Parameters** `mentions` – A list of strings

:return True/False

**prune\_by\_score** (*query\_inss*)

prunes based on a static threshold of ranking score.

**prune\_containment\_mentions** (*query\_inss*)

Deletes containment mentions, if they have lower score.

## nordlys.logic.el.ltr module

### LTR Entity Linking Approach

Class for Learning-to-Rank entity linking approach

**Author** Faegheh Hasibi

**class** `nordlys.logic.el.ltr.LTR` (*query, entity, elastic, fcache, model=None, threshold=None, cmns\_th=0.1*)

Bases: `object`

**disambiguate** (*inss*)

Performs disambiguation

**static gen\_train\_set** (*gt, query\_file, train\_set*)

Trains LTR model for entity linking.

**get\_candidate\_inss** ()

Detects mentions and their candidate entities (with their commonness scores) and generates instances

**Returns** Instances object

**get\_features** (*ins, cand\_ens=None*)

Generates the features set for each instance.

**Parameters**

- **ins** – instance object
- **cand\_ens** – dictionary of candidate entities {en\_id: cmns, ... }

**Returns** dictionary of features {ftr\_name: value, ... }

**link** ()

Links the query to the entity.

**Returns** dictionary [{"mention": xx, "entity": yy, "score": zz}, ... ]

**static load\_yerd** (*gt\_file*)

Reads the Y-ERD collection and returns a dictionary.

**Parameters** **gt\_file** – Path to the Y-ERD collection

**Returns** dictionary {(qid, query, en\_id, mention) ... }

**rank\_ens** ()

Ranks instances according to the learned LTR model

**Parameters** **n** – length of n-gram

**Returns** dictionary {(dbp\_uri, fb\_id):commonness, ..}

**static train** (*config*)

## nordlys.logic.elr package

### ELR

This package is the implementation of ELR-based models.

### Submodules

#### nordlys.logic.elr.field\_mapping module

### Field Mapping for ELR

Computes PRMS field mapping probabilities.

**Author** Faegheh Hasibi

```
class nordlys.logic.elr.field_mapping.FieldMapping (elastic_uri, n)
    Bases: object

    DEBUG = 0
    MAPPING_DEBUG = 0

    map (en_id)
        Gets PRMS mapping probability for a clique type
            Return Dictionary {field weight, ..}

nordlys.logic.elr.field_mapping.arg_parser()
nordlys.logic.elr.field_mapping.load_entities (annot_file, th=0.1)
nordlys.logic.elr.field_mapping.main (args)
```

**nordlys.logic.elr.scorer\_elr** module

**nordlys.logic.elr.top\_fields** module

## Top Fields

This class returns top fields based on document frequency

**Author** Faegheh Hasibi

```
class nordlys.logic.elr.top_fields.TopFields (elastic)
    Bases: object

    DEBUG = 0
    fields

    get_top_term (en, n)
        Returns top-n fields with highest document frequency for the given entity ID.
```

**nordlys.logic.entity** package

## Entity

This is the entity package.

## Submodules

**nordlys.logic.entity.entity** module

## Entity

Provides access to entity catalogs (DBpedia and surface forms).

**Author** Faegheh Hasibi

```
class nordlys.logic.entity.entity.Entity
    Bases: object
```

**dbp\_to\_fb** (*dbp\_id*)

Converts DBpedia id to Freebase; it returns list of Freebase IDs.

**fb\_to\_dbp** (*fb\_id*)

Converts Freebase id to DBpedia; it returns list of DBpedia IDs.

**lookup\_en** (*entity\_id*)

Looks up an entity by its identifier.

**Parameters** *entity\_id* – entity identifier (“<dbpedia:Audi\_A4>”)

:return A dictionary with the entity document or None.

**lookup\_name\_dbpedia** (*name*)

Looks up a name in a surface form dictionary and returns all candidate entities.

**lookup\_name\_facc** (*name*)

Looks up a name in a surface form dictionary and returns all candidate entities.

## nordlys.logic.er package

### Entity retrieval

This is the entity retrieval package.

### Submodules

#### nordlys.logic.er.entity\_retrieval module

#### nordlys.logic.er.field\_mapping module

### Field Mapping for ER

Computes PRMS field mapping probabilities.

**Author** Faegheh Hasibi

**class** nordlys.logic.er.field\_mapping.**FieldMapping** (*elastic\_uri, n*)

Bases: `object`

**DEBUG** = 0

**MAPPING\_DEBUG** = 0

**map** (*en\_id*)

Gets PRMS mapping probability for a clique type

**Return Dictionary** {*field* *weight*, ..}

nordlys.logic.er.field\_mapping.**arg\_parser** ()

nordlys.logic.er.field\_mapping.**load\_entities** (*annot\_file, th=0.1*)

nordlys.logic.er.field\_mapping.**main** (*args*)

## nordlys.logic.er.scorer\_elr module

## nordlys.logic.er.top\_fields module

### Top Fields

This class returns top fields based on document frequency

**Author** Faegheh Hasibi

**class** nordlys.logic.er.top\_fields.**TopFields** (*elastic*)

Bases: `object`

**DEBUG** = 0

**fields**

**get\_top\_term** (*en, n*)

Returns top-n fields with highest document frequency for the given entity ID.

## nordlys.logic.features package

### Features

This is the features package.

### Submodules

## nordlys.logic.features.feature\_cache module

### Feature

Implements a generic feature class.

Authors: Faegheh Hasibi

**class** nordlys.logic.features.feature\_cache.**FeatureCache**

Bases: `object`

**get\_feature\_val** (*feature\_name, key, callback\_func, \*args*)

Checks the cache and computes the feature if it does not exist

**set\_feature\_val** (*feature\_name, key, value*)

Adds a feature and its value to the cache.

#### Parameters

- **feature\_name** – Name of the feature
- **key** – the name of what feature is computed for (e.g., a mention, entity)
- **value** – feature value

## nordlys.logic.features.ftr\_entity module

### FTR Entity

Implements features related to an entity.

**Author** Faegheh Hasibi

**class** nordlys.logic.features.ftr\_entity.**FtrEntity**(*en\_id, entity*)

Bases: `object`

**outlinks** ()

Number of entity out-links

**redirects** ()

Number of redirect pages linking to the entity

## nordlys.logic.features.ftr\_entity\_mention module

### FTR Entity Mention

Implements features related to an entity-mention pair.

**Author** Faegheh Hasibi

**class** nordlys.logic.features.ftr\_entity\_mention.**FtrEntityMention**(*en\_id, mention, entity*)

Bases: `object`

**commonness** ()

Computes probability of entity *e* being linked by mention:  $\text{link}(e,m)/\text{link}(m)$  Returns zero if  $\text{link}(m) = 0$

**mct** ()

Returns True if mention contains the title of entity

**pos1** ()

Returns position of the occurrence of mention in the short abstract.

**tcm** ()

Returns True if title of entity contains mention

**tem** ()

Returns True if title of entity equals mention.

## nordlys.logic.features.ftr\_entity\_similarity module

### FTR Entity Similarity

Implements features capturing the similarity between entity and a query.

**Author** Faegheh Hasibi

**class** nordlys.logic.features.ftr\_entity\_similarity.**FtrEntitySimilarity**(*query, en\_id, elastic*)

Bases: `object`

`DEBUG = 0`

`context_sim(mention, field='catchall')`

**LM score of entity to the context of query (context means query - mention)**

E.g. given the query “uss yorktown charleston” and mention “uss”, query context is ” yorktown charleston”

**Parameters**

- **mention** – string
- **field** – field name

:return context similarity score

`lm_score(field='catchall')`

Query length normalized LM score between entity field and query

**Parameters** **field** – field name

:return MLM score

`mlm_score(field_weights)`

Query length normalized MLM similarity between the entity and query

**Parameters** **field\_weights** – dictionary {field: weight, ... }

:return MLM score

`nllr(query, field_weights)`

**Computes Normalized query likelihood (NLLR):**  $NLLR(q,d) = \sum_{\{t \text{ in } q\}} P(t|q) \log P(t|C) - \sum_{\{t \text{ in } q\}} p(t|q) \log P(t|C)$  where:

$$P(t|q) = n(t,q)/|q| \quad P(t|C) = \sum_{\{f\}} \mu_f * P(t|C_f) \quad P(t|C) = \text{smoothed LM/MLM score}$$

**Parameters**

- **query** – query
- **field\_weights** – dictionary {field: weight, ... }

**Returns** NLLR score

**nordlys.logic.features.ftr\_lexical module**

**nordlys.logic.features.ftr\_mention module**

**FTR Mention**

Implements mention feature.

**Author** Faegheh Hasibi

**class** nordlys.logic.features.ftr\_mention.**FtrMention**(*mention*, *entity=None*, *cand\_ens=None*)

Bases: `object`

**len\_ratio**(*q*)

Computes mention to query length.

**matches** ()

Number of entities whose surface form equals the mention. Uses both DBpedia and Freebase name variants.

**mention\_len** ()

Number of terms in the mention

## nordlys.logic.features.word2vec module

### Word2vec

Implements functionalities over the 300-dim GoogleNews word2vec semantic representations of words.

**Author** Dario Garigliotti

**class** nordlys.logic.features.word2vec.**Word2Vec** (*mongo*)

Bases: `object`

**get\_centroid\_vector** (*s*)

Returns the normalized sum of the word2vec vectors corresponding to the terms in *s*.

**Parameters** *s* (*str*) – a phrase.

**Returns** Centroid vector of the terms in *s*.

**get\_vector** (*word*)

Gets the w2v vector corresponding to the word, or a zero-valued vector if not present.

**Parameters** *word* (*str*) – a word.

**Returns**

nordlys.logic.features.word2vec.**arg\_parser** ()

nordlys.logic.features.word2vec.**main** (*args*)

## nordlys.logic.fusion package

### Submodules

#### nordlys.logic.fusion.fusion\_scorer module

### Fusion Scorer

Abstract class for fusion-based scoring.

**Authors** Shuo Zhang, Krisztian Balog, Dario Garigliotti

**class** nordlys.logic.fusion.fusion\_scorer.**FusionScorer** (*index\_name*, *association\_file=*`None`, *run\_id=*`'fusion'`)

Bases: `object`

**Parameters**

- **index\_name** – name of index
- **association\_file** – association file

**ASSOC\_MODE\_BINARY** = 1

**ASSOC\_MODE\_UNIFORM** = 2

Abstract class for any fusion-based method.

**load\_associations** ()

Loads the document-object associations.

**load\_queries** (*query\_file*)

Loads the query file :return: query dictionary {queryID:query([term1,term2,...])}

**score\_queries** (*queries, output\_file*)

Scores all queries and optionally dumps results into an output file.

**score\_query** (*query, assoc\_fun=None*)

## nordlys.logic.fusion.late\_fusion\_scorer module

### Late Fusion Scorer

Class for late fusion scorer (i.e., document-centric model).

**Authors** Shuo Zhang, Krisztian Balog, Dario Garigliotti

```
class nordlys.logic.fusion.late_fusion_scorer.LateFusionScorer (index_name,  
retr_model,  
retr_params,  
num_docs=None,  
field='content',  
run_id='fusion',  
num_objs=100,  
assoc_mode=1,  
as-  
soc_file=None)
```

Bases: *nordlys.logic.fusion.fusion\_scorer.FusionScorer*

#### Parameters

- **index\_name** – name of index
- **assoc\_file** – document-object association file
- **assoc\_mode** – document-object weight mode, uniform or binary
- **retr\_model** – the retrieval model; valid values: “lm”, “bm25”
- **retr\_params** – config including smoothing method and parameter
- **num\_objs** – the number of ranked objects for a query
- **assoc\_mode** – the fusion weights, which could be binary or uniform
- **assoc\_file** – object-doc association file

**score\_query** (*query, assoc\_fun=None*)

Scores a given query.

**Parameters** **query** – query string.

**Returns** a RetrievalResults instance.

**Func** **assoc\_fun** function to return a list of docs for an object

## nordlys.logic.query package

### Query

This is the query package.

### Submodules

#### nordlys.logic.query.mention module

### Mention

Class for entity mentions (used for entity linking)

- Generates all candidate entities for a mention
- Computes commonness for a mention-entity pairs

**class** nordlys.logic.query.mention.**Mention** (*mention, entity, cmns\_th=None*)

Bases: `object`

**get\_cand\_ens** ()

Returns all candidate entities for the mention

**Returns** {en:cmn\_score}

nordlys.logic.query.mention.**main** (*args*)

#### nordlys.logic.query.query module

### Query

Class for representing a query.

TODO: add preprocessing using

**Author** Faegheh Hasibi

**class** nordlys.logic.query.query.**Query** (*query, qid=""*)

Bases: `object`

**get\_ngrams** ()

Finds all n-grams of the query.

**Returns** list of n-grams

**get\_terms** ()

Gets query terms.

**Returns** list of query terms

**qid**

**query**

**raw\_query**

## nordlys.logic.tti package

### Submodules

#### nordlys.logic.tti.type\_centric module

Type centric method for TTI.

@author:

```
class nordlys.logic.tti.type_centric.TypeCentric (query, retrieval_config)  
    Bases: object
```

## nordlys.services package

### Services

All modules in this package can be used from the command line or from a RESTful API.

### Submodules

#### nordlys.services.api module

#### Nordlys API

This is the main console application for the Nordlys API.

**Authors** Krisztian Balog, Faegheh Hasibi, Shuo Zhang

nordlys.services.api.**after\_request** (*response*)

nordlys.services.api.**catalog\_dbp2fb** (*dbp\_id*)

nordlys.services.api.**catalog\_fb2dbp** (*fb\_id*)

nordlys.services.api.**catalog\_lookup\_id** (*entity\_id*)

nordlys.services.api.**catalog\_lookup\_sf\_dbpedia** (*sf*)

nordlys.services.api.**catalog\_lookup\_sf\_facc** (*sf*)

nordlys.services.api.**entity\_linking** ()

nordlys.services.api.**entity\_types** ()

nordlys.services.api.**error** (*str*)

@todo complete error handling

**Parameters** *str* –

**Returns**

nordlys.services.api.**exceptions** (*e*)

nordlys.services.api.**index** ()

nordlys.services.api.**retrieval** ()

## nordlys.services.ec module

### Entity catalog

Command line end point for entity catalog

### Usage

```
python -m nordlys.services.ec -o <operation> -i <input>
```

### Examples

- `python -m nordlys.services.ec -o lookup_id -i <dbpedia:Audi_A4>`
- `python -m nordlys.services.ec -o "lookup_sf_dbpedia" -i "audi a4"`
- `python -m nordlys.services.ec -o "lookup_sf_facc" -i "audi a4"`
- `python -m nordlys.services.ec -o "dbpedia2freebase" -i "<dbpedia:Audi_A4>"`
- `python -m nordlys.services.ec -o "freebase2dbpedia" -i "<fb:m.030qmx>"`

**Author** Faegheh Hasibi

```
nordlys.services.ec.arg_parser ()
```

```
nordlys.services.ec.main (args)
```

## nordlys.services.el module

### Entity Linking

The command-line application for entity linking

### Usage

```
python -m nordlys.services.el -c <config_file> -q <query>
```

If `-q <query>` is passed, it returns the results for the specified query and prints them in terminal.

### Config parameters

- **method: name of the method**
  - **cmns** The baseline method that uses the overall popularity of entities as link targets
  - **ltr** The learning-to-rank model
- **threshold:** Entity linking threshold; varies depending on the method (*default: 0.1*)
- **step:** The step of entity linking process: [linking|ranking|disambiguation], (*default: linking*)
- **kb\_snapshot:** File containing the KB snapshot of proper named entities; required for LTR, and optional for CMNS

- **query\_file**: name of query file (JSON)
- **output\_file**: name of output file

Parameters of LTR method:

- **model\_file**: The trained model file; (*default: "data/el/model.txt"*)
- **ground\_truth**: The ground truth file; (*optional*)
- **gen\_training\_set**: If True, generates the training set from the groundtruth and query files; (*default: False*)
- **gen\_model**: If True, trains the model from the training set; (*default: False*)
- The other parameters are similar to the nordlys.core.ml.ml settings

### Example config

```
{
  "method": "cmns",
  "threshold": 0.1,
  "query_file": "path/to/queries.json"
  "output_file": "path/to/output.json"
}
```

---

**Author** Faegheh Hasibi

**class** nordlys.services.el.**EL** (*config, entity, elastic=None, fcache=None*)

Bases: `object`

**batch\_linking** ()

Scores queries in a batch and outputs results.

**link** (*query, qid=""*)

Performs entity linking for the query.

**Parameters** **query** – query string

**Returns** annotated query

nordlys.services.el.**arg\_parser** ()

nordlys.services.el.**main** (*args*)

### nordlys.services.er module

#### Entity Retrieval

Command-line application for entity retrieval.

#### Usage

```
python -m nordlys.services.er -c <config_file> -q <query>
```

If `-q <query>` is passed, it returns the results for the specified query and prints them in terminal.

## Config parameters

- **index\_name**: name of the index,
- **first\_pass**:
  - **num\_docs**: number of documents in first-pass scoring (default: 100)
  - **field**: field used in first pass retrieval (default: Elastic.FIELD\_CATCHALL)
  - **fields\_return**: comma-separated list of fields to return for each hit (default: “”)
- **num\_docs**: number of documents to return (default: 100)
- **start**: starting offset for ranked documents (default:0)
- **model**: name of retrieval model; accepted values: [lm, mlm, prms] (default: lm)
- **field**: field name for LM (default: catchall)
- **fields**: list of fields for PRMS (default: [catchall])
- **field\_weights**: dictionary with fields and corresponding weights for MLM (default: {catchall: 1})
- **smoothing\_method**: accepted values: [jm, dirichlet] (default: dirichlet)
- **smoothing\_param**: value of lambda or mu; accepted values: [float or “avg\_len”], (jm default: 0.1, dirichlet default: 2000)
- **query\_file**: name of query file (JSON),
- **output\_file**: name of output file,
- **run\_id**: run id for TREC output

## Example config

```
{
  "index_name": "dbpedia_2015_10",
  "first_pass": {
    "num_docs": 1000
  },
  "model": "prms",
  "num_docs": 1000,
  "smoothing_method": "dirichlet",
  "smoothing_param": 2000,
  "fields": ["names", "categories", "attributes", "similar_entity_names", "related_
↪entity_names"],
  "query_file": "path/to/queries.json",
  "output_file": "path/to/output.txt",
  "run_id": "test"
}
```

**Author** Faegheh Hasibi

```
class nordlys.services.er.ER(config, elastic=None)
    Bases: object

    batch_retrieval()
        Performs batch retrieval for a set of queries
```

**retrieve** (*query*)

Retrieves entities for a query

```
nordlys.services.er.arg_parser()
```

```
nordlys.services.er.main(args)
```

### nordlys.services.tti module

## Target Type Identification

The command-line application for target type identification.

### Usage

```
python -m nordlys.services.tti <config_file> -q <query>
```

If `-q <query>` is passed, it returns the results for the specified query and prints them in terminal.

### Config parameters

- **method**: name of TTI method; accepted values: ["tc", "ec", "ltr"]
- **num\_docs**: number of documents to return
- **start**: starting offset for ranked documents
- **model**: retrieval model, if method is "tc" or "ec"; accepted values: ["lm", "bm25"]
- **ec\_cutoff**: if method is "ec", rank cut-off of top-*K* entities for EC TTI
- **field**: field name, if method is "tc" or "ec"
- **smoothing\_method**: accepted values: ["jm", "dirichlet"]
- **smoothing\_param**: value of lambda or mu; accepted values: [float or "avg\_len"]
- **query\_file**: path to query file (JSON)
- **output\_file**: path to output file (JSON)
- **trec\_output\_file**: path to output file (trec\_eval-formatted)

### Example config

```
{ "method": "ec",  
  "num_docs": 10,  
  "model": "lm",  
  "first_pass": {  
    "num_docs": 50  
  },  
  "smoothing_method": "dirichlet",  
  "smoothing_param": 2000,  
  "ec_cutoff": 20,  
  "query_file": "path/to/queries.json",
```

(continues on next page)

(continued from previous page)

```
"output_file": "path/to/output.txt",  
}
```

**Author** Dario Garigliotti

**class** nordlys.services.tti.**TTI** (*config*)

Bases: `object`

**batch\_identification** ()

Annotates, in a batch, queries with identified target types, and outputs results.

**identify** (*query*)

Performs target type identification for the query.

**Parameters** **query** (*str*) – query string

**Returns** annotated query

nordlys.services.tti.**arg\_parser** ()

nordlys.services.tti.**main** (*args*)

## 2.7.2 Submodules

### nordlys.config module

#### config

Global nordlys config.

**Author** Krisztian Balog

**Author** Faegheh Hasibi

nordlys.config.**load\_nordlys\_config** (*file\_name*)

Loads nordlys config file. If local file is provided, global one is ignored.

## 2.8 Contact

Nordlys is written by:

- Faegheh Hasibi
- Krisztian Balog
- Dario Garigliotti
- Shuo Zhang

For any queries contact <faegheh.hasibi@ntnu.no> or <krisztian.balog@uis.no>.

*Drop us some words! We are looking forward to hear your feedback!*



## CHAPTER 3

---

### Indices and tables

---

- genindex
- modindex
- search



**n**

nordlys, 22  
 nordlys.config, 69  
 nordlys.core, 22  
 nordlys.core.data, 22  
 nordlys.core.data.dbpedia, 22  
 nordlys.core.data.dbpedia.indexer\_dbpedia\_types, 23  
 nordlys.core.data.facc, 23  
 nordlys.core.data.facc.facc2mongo, 23  
 nordlys.core.data.word2vec, 24  
 nordlys.core.data.word2vec.word2vec2mongo, 24  
 nordlys.core.eval, 24  
 nordlys.core.eval.eval, 25  
 nordlys.core.eval.plot\_diff, 25  
 nordlys.core.eval.query\_diff, 26  
 nordlys.core.eval.trec\_eval, 26  
 nordlys.core.eval.trec\_qrels, 27  
 nordlys.core.eval.trec\_run, 28  
 nordlys.core.ml, 29  
 nordlys.core.ml.cross\_validation, 30  
 nordlys.core.ml.instance, 31  
 nordlys.core.ml.instances, 33  
 nordlys.core.ml.ml, 35  
 nordlys.core.retrieval, 37  
 nordlys.core.retrieval.elastic, 38  
 nordlys.core.retrieval.elastic\_cache, 42  
 nordlys.core.retrieval.indexer\_mongo, 43  
 nordlys.core.retrieval.retrieval, 43  
 nordlys.core.retrieval.retrieval\_results, 45  
 nordlys.core.retrieval.scorer, 45  
 nordlys.core.retrieval.toy\_indexer, 48  
 nordlys.core.storage, 48  
 nordlys.core.storage.mongo, 50  
 nordlys.core.storage.parser, 48  
 nordlys.core.storage.parser.nt\_parser, 48  
 nordlys.core.storage.parser.uri\_prefix, 49  
 nordlys.core.utils, 51  
 nordlys.core.utils.file\_utils, 51  
 nordlys.core.utils.logging\_utils, 52  
 nordlys.logic, 53  
 nordlys.logic.el, 53  
 nordlys.logic.el.cmns, 53  
 nordlys.logic.el.el\_utils, 53  
 nordlys.logic.el.greedy, 54  
 nordlys.logic.el.ltr, 54  
 nordlys.logic.elr, 55  
 nordlys.logic.elr.field\_mapping, 55  
 nordlys.logic.elr.top\_fields, 56  
 nordlys.logic.entity, 56  
 nordlys.logic.entity.entity, 56  
 nordlys.logic.er, 57  
 nordlys.logic.er.field\_mapping, 57  
 nordlys.logic.er.top\_fields, 58  
 nordlys.logic.features, 58  
 nordlys.logic.features.feature\_cache, 58  
 nordlys.logic.features.ftr\_entity, 59  
 nordlys.logic.features.ftr\_entity\_mention, 59  
 nordlys.logic.features.ftr\_entity\_similarity, 59  
 nordlys.logic.features.ftr\_mention, 60  
 nordlys.logic.features.word2vec, 61  
 nordlys.logic.fusion, 61  
 nordlys.logic.fusion.fusion\_scorer, 61  
 nordlys.logic.fusion.late\_fusion\_scorer, 62  
 nordlys.logic.query, 63  
 nordlys.logic.query.mention, 63  
 nordlys.logic.query.query, 63  
 nordlys.logic.tti, 64  
 nordlys.logic.tti.type\_centric, 64

nordlys.services, 64  
nordlys.services.api, 64  
nordlys.services.ec, 65  
nordlys.services.el, 65  
nordlys.services.er, 66  
nordlys.services.tti, 68

## A

- add() (*nordlys.core.storage.mongo.Mongo* method), 50  
 add\_doc() (*nordlys.core.retrieval.elastic.Elastic* method), 39  
 add\_docs\_bulk() (*nordlys.core.retrieval.elastic.Elastic* method), 39  
 add\_feature() (*nordlys.core.ml.instance.Instance* method), 32  
 add\_features\_from\_tsv() (*nordlys.core.ml.instances.Instances* method), 33  
 add\_instance() (*nordlys.core.ml.instances.Instances* method), 33  
 add\_properties\_from\_tsv() (*nordlys.core.ml.instances.Instances* method), 33  
 add\_property() (*nordlys.core.ml.instance.Instance* method), 32  
 add\_qids() (*nordlys.core.ml.instances.Instances* method), 33  
 add\_target\_from\_tsv() (*nordlys.core.ml.instances.Instances* method), 33  
 after\_request() (in module *nordlys.services.api*), 64  
 analyse\_features() (*nordlys.core.ml.ml.ML* method), 36  
 analyze\_query() (*nordlys.core.retrieval.elastic.Elastic* method), 39  
 analyzed\_field() (*nordlys.core.retrieval.elastic.Elastic* static method), 39  
 ANALYZER\_STOP (*nordlys.core.retrieval.elastic.Elastic* attribute), 39  
 ANALYZER\_STOP\_STEM (*nordlys.core.retrieval.elastic.Elastic* attribute), 39  
 append() (*nordlys.core.retrieval.retrieval\_results.RetrievalResults* method), 45  
 append\_dict() (*nordlys.core.storage.mongo.Mongo* method), 50  
 append\_instances() (*nordlys.core.ml.instances.Instances* method), 33  
 append\_list() (*nordlys.core.storage.mongo.Mongo* method), 50  
 append\_set() (*nordlys.core.storage.mongo.Mongo* method), 50  
 apply\_model() (*nordlys.core.ml.ml.ML* method), 36  
 arg\_parser() (in module *nordlys.core.data.dbpedia.indexer\_dbpedia\_types*), 23  
 arg\_parser() (in module *nordlys.core.data.facc.facc2mongo*), 24  
 arg\_parser() (in module *nordlys.core.data.word2vec.word2vec2mongo*), 24  
 arg\_parser() (in module *nordlys.core.eval.eval*), 25  
 arg\_parser() (in module *nordlys.core.eval.trec\_qrels*), 28  
 arg\_parser() (in module *nordlys.core.eval.trec\_run*), 29  
 arg\_parser() (in module *nordlys.core.ml.ml*), 37  
 arg\_parser() (in module *nordlys.core.retrieval.retrieval*), 45  
 arg\_parser() (in module *nordlys.logic.elr.field\_mapping*), 56  
 arg\_parser() (in module *nordlys.logic.er.field\_mapping*), 57  
 arg\_parser() (in module *nordlys.logic.features.word2vec*), 61  
 arg\_parser() (in module *nordlys.services.ec*), 65  
 arg\_parser() (in module *nordlys.services.el*), 66  
 arg\_parser() (in module *nordlys.services.er*), 68  
 arg\_parser() (in module *nordlys.services.tti*), 69  
 ASSOC\_MODE\_BINARY (*nordlys.logic.fusion.fusion\_scorer.FusionScorer* attribute), 61  
 ASSOC\_MODE\_UNIFORM (*nordlys.logic.fusion.fusion\_scorer.FusionScorer*

- attribute*), 62
  - avg\_len() (*nordlys.core.retrieval.elastic.Elastic method*), 39
  - avg\_len() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
- B**
- batch\_identification() (*nordlys.services.tti.TTI method*), 69
  - batch\_linking() (*nordlys.services.el.EL method*), 66
  - batch\_retrieval() (*nordlys.core.retrieval.retrieval.Retrieval method*), 44
  - batch\_retrieval() (*nordlys.services.er.ER method*), 67
  - BM25 (*nordlys.core.retrieval.elastic.Elastic attribute*), 39
  - build() (*nordlys.core.data.facc.facc2mongo.FACCToMongo method*), 24
  - build() (*nordlys.core.data.word2vec.word2vec2mongo.Word2VecToMongo method*), 24
  - build() (*nordlys.core.retrieval.indexer\_mongo.IndexerMongo method*), 43
  - build\_index() (*nordlys.core.data.dbpedia.indexer\_dbpedia.IndexerDbpedia method*), 23
- C**
- catalog\_dbp2fb() (*in module nordlys.services.api*), 64
  - catalog\_fb2dbp() (*in module nordlys.services.api*), 64
  - catalog\_lookup\_id() (*in module nordlys.services.api*), 64
  - catalog\_lookup\_sf\_dbpedia() (*in module nordlys.services.api*), 64
  - catalog\_lookup\_sf\_facc() (*in module nordlys.services.api*), 64
  - check\_config() (*nordlys.core.retrieval.retrieval.Retrieval static method*), 44
  - Cmns (*class in nordlys.logic.el.cmns*), 53
  - coll\_length() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - coll\_length() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
  - coll\_term\_freq() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - coll\_term\_freq() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
  - commonness() (*nordlys.logic.features.ftr\_entity\_mention.FtrEntityMention method*), 59
  - context\_sim() (*nordlys.logic.features.ftr\_entity\_similarity.FtrEntitySimilarity method*), 60
  - convert\_txt\_to\_json() (*in module nordlys.core.storage.parser.uri\_prefix*), 50
- create\_folds() (*nordlys.core.ml.cross\_validation.CrossValidation method*), 30
  - create\_index() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - create\_interpretations() (*nordlys.logic.el.greedy.Greedy method*), 54
  - create\_pdf() (*nordlys.core.eval.plot\_diff.QueryDiff method*), 25
  - CrossValidation (*class in nordlys.core.ml.cross\_validation*), 30
- D**
- dbp\_to\_fb() (*nordlys.logic.entity.entity.Entity method*), 56
  - DEBUG (*nordlys.logic.elr.field\_mapping.FieldMapping attribute*), 56
  - DEBUG (*nordlys.logic.elr.top\_fields.TopFields attribute*), 56
  - DEBUG (*nordlys.logic.er.field\_mapping.FieldMapping attribute*), 57
  - DEBUG (*nordlys.logic.er.top\_fields.TopFields attribute*), 58
  - DEBUG (*nordlys.logic.elr.top\_fields.TopFields attribute*), 59
  - delete\_index() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - DIRICHLET (*nordlys.core.retrieval.scorer.ScorerLM attribute*), 46
  - disambiguate() (*nordlys.logic.el.cmns.Cmns method*), 53
  - disambiguate() (*nordlys.logic.el.greedy.Greedy method*), 54
  - disambiguate() (*nordlys.logic.el.ltr.LTR method*), 54
  - doc\_count() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - doc\_count() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
  - doc\_freq() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - doc\_freq() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
  - doc\_length() (*nordlys.core.retrieval.elastic.Elastic method*), 40
  - doc\_length() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
  - DOC\_TYPE (*nordlys.core.retrieval.elastic.Elastic attribute*), 39
  - drop() (*nordlys.core.storage.mongo.Mongo method*), 51
  - dump\_differences() (*nordlys.core.eval.query\_diff.QueryDiff method*), 26

`dump_tsv()` (*nordlys.core.utils.file\_utils.FileUtils static method*), 52

## E

`EL` (*class in nordlys.services.el*), 66

`Elastic` (*class in nordlys.core.retrieval.elastic*), 39

`elastic_to_retrieval()`

(*nordlys.core.retrieval.retrieval\_results.RetrievalResults class method*), 45

`ElasticCache` (*class in nordlys.core.retrieval.elastic\_cache*), 42

`Entity` (*class in nordlys.logic.entity.entity*), 56

`entity_linking()` (*in module nordlys.services.api*), 64

`entity_types()` (*in module nordlys.services.api*), 64

`ER` (*class in nordlys.services.er*), 67

`error()` (*in module nordlys.services.api*), 64

`Eval` (*class in nordlys.core.eval.eval*), 25

`evaluate()` (*nordlys.core.eval.trec\_eval.TrecEval method*), 26

`exceptions()` (*in module nordlys.services.api*), 64

## F

`FACCToMongo` (*class in nordlys.core.data.facc.facc2mongo*), 24

`fb_to_dbp()` (*nordlys.logic.entity.entity.Entity method*), 57

`FeatureCache` (*class in nordlys.logic.features.feature\_cache*), 58

`features` (*nordlys.core.ml.instance.Instance attribute*), 32

`FIELD_CATCHALL` (*nordlys.core.retrieval.elastic.Elastic attribute*), 39

`FIELD_ELASTIC_CATCHALL` (*nordlys.core.retrieval.elastic.Elastic attribute*), 39

`FIELDDED_MODELS` (*nordlys.core.retrieval.retrieval.Retrieval attribute*), 44

`FieldMapping` (*class in nordlys.logic.elr.field\_mapping*), 55

`FieldMapping` (*class in nordlys.logic.er.field\_mapping*), 57

`fields` (*nordlys.logic.elr.top\_fields.TopFields attribute*), 56

`fields` (*nordlys.logic.er.top\_fields.TopFields attribute*), 58

`FileUtils` (*class in nordlys.core.utils.file\_utils*), 52

`filter()` (*nordlys.core.eval.trec\_run.TrecRun method*), 28

`filter_by_doc_ids()` (*nordlys.core.eval.trec\_qrels.TrecQrels method*), 27

`filter_by_query_ids()` (*nordlys.core.eval.trec\_qrels.TrecQrels*

*method*), 27

`find_all()` (*nordlys.core.storage.mongo.Mongo method*), 51

`find_by_id()` (*nordlys.core.storage.mongo.Mongo method*), 51

`from_json()` (*nordlys.core.ml.instance.Instance class method*), 32

`from_json()` (*nordlys.core.ml.instances.Instances class method*), 33

`FtrEntity` (*class in nordlys.logic.features.ftr\_entity*), 59

`FtrEntityMention` (*class in nordlys.logic.features.ftr\_entity\_mention*), 59

`FtrEntitySimilarity` (*class in nordlys.logic.features.ftr\_entity\_similarity*), 59

`FtrMention` (*class in nordlys.logic.features.ftr\_mention*), 60

`FusionScorer` (*class in nordlys.logic.fusion.fusion\_scorer*), 61

## G

`gen_model()` (*nordlys.core.ml.ml.ML method*), 36

`gen_train_set()` (*nordlys.logic.el.ltr.LTR static method*), 55

`get_all()` (*nordlys.core.ml.instances.Instances method*), 33

`get_all_ids()` (*nordlys.core.ml.instances.Instances method*), 34

`get_cand_ens()` (*nordlys.logic.query.mention.Mention method*), 63

`get_candidate_inss()` (*nordlys.logic.el.ltr.LTR method*), 55

`get_centroid_vector()` (*nordlys.logic.features.word2vec.Word2Vec method*), 61

`get_config()` (*in module nordlys.core.retrieval.retrieval*), 45

`get_dirichlet_prob()` (*nordlys.core.retrieval.scorer.ScorerLM static method*), 46

`get_doc()` (*nordlys.core.retrieval.elastic.Elastic method*), 40

`get_feature()` (*nordlys.core.ml.instance.Instance method*), 32

`get_feature_val()` (*nordlys.logic.features.feature\_cache.FeatureCache method*), 58

`get_features()` (*nordlys.logic.el.ltr.LTR method*), 55

`get_field_stats()` (*nordlys.core.retrieval.elastic.Elastic method*), 40

get\_fields() (*nordlys.core.retrieval.elastic.Elastic method*), 40  
 get\_folds() (*nordlys.core.ml.cross\_validation.CrossValidation method*), 31  
 get\_instance() (*nordlys.core.ml.instances.Instances method*), 34  
 get\_instances() (*nordlys.core.ml.cross\_validation.CrossValidation method*), 31  
 get\_jm\_prob() (*nordlys.core.retrieval.scorer.ScorerLM static method*), 46  
 get\_lm\_term\_prob() (*nordlys.core.retrieval.scorer.ScorerLM method*), 46  
 get\_lm\_term\_probs() (*nordlys.core.retrieval.scorer.ScorerLM method*), 46  
 get\_mapping() (*nordlys.core.retrieval.elastic.Elastic method*), 40  
 get\_mapping\_prob() (*nordlys.core.retrieval.scorer.ScorerPRMS method*), 47  
 get\_mapping\_probs() (*nordlys.core.retrieval.scorer.ScorerPRMS method*), 47  
 get\_mlm\_term\_prob() (*nordlys.core.retrieval.scorer.ScorerMLM method*), 47  
 get\_mlm\_term\_probs() (*nordlys.core.retrieval.scorer.ScorerMLM method*), 47  
 get\_ngrams() (*nordlys.logic.query.query.Query method*), 63  
 get\_num\_docs() (*nordlys.core.storage.mongo.Mongo method*), 51  
 get\_prefixed() (*nordlys.core.storage.parser.uri\_prefix.UriPrefix method*), 50  
 get\_property() (*nordlys.core.ml.instance.Instance method*), 32  
 get\_queries() (*nordlys.core.eval.trec\_qrels.TrecQrels method*), 27  
 get\_query\_ids() (*nordlys.core.eval.trec\_eval.TrecEval method*), 26  
 get\_query\_results() (*nordlys.core.eval.trec\_run.TrecRun method*), 29  
 get\_rel() (*nordlys.core.eval.trec\_qrels.TrecQrels method*), 27  
 get\_results() (*nordlys.core.eval.trec\_run.TrecRun method*), 29  
 get\_score() (*nordlys.core.eval.trec\_eval.TrecEval method*), 26  
 get\_score() (*nordlys.core.retrieval.retrieval\_results.RetrievalResults method*), 45  
 get\_scorer() (*nordlys.core.retrieval.scorer.Scorer static method*), 45  
 get\_scores\_sorted() (*nordlys.core.retrieval.retrieval\_results.RetrievalResults method*), 45  
 get\_settings() (*nordlys.core.retrieval.elastic.Elastic method*), 40  
 get\_top\_term() (*nordlys.logic.query.query.Query method*), 63  
 get\_top\_term() (*nordlys.logic.elr.top\_fields.TopFields method*), 56  
 get\_top\_term() (*nordlys.logic.er.top\_fields.TopFields method*), 58  
 get\_total\_field\_freq() (*nordlys.core.retrieval.scorer.ScorerPRMS method*), 47  
 get\_vector() (*nordlys.logic.features.word2vec.Word2Vec method*), 61  
 Greedy (*class in nordlys.logic.el.greedy*), 54  
 group\_by\_property() (*nordlys.core.ml.instances.Instances method*), 34

**I**

id (*nordlys.core.ml.instance.Instance attribute*), 32  
 ID\_FIELD (*nordlys.core.storage.mongo.Mongo attribute*), 50  
 identify() (*nordlys.services.tti.TTI method*), 69  
 inc() (*nordlys.core.storage.mongo.Mongo method*), 51  
 inc\_in\_dict() (*nordlys.core.storage.mongo.Mongo method*), 51  
 index() (*in module nordlys.services.api*), 64  
 IndexerDBpediaTypes (*class in nordlys.core.data.dbpedia.indexer\_dbpedia\_types*), 23  
 IndexerMongo (*class in nordlys.core.retrieval.indexer\_mongo*), 43  
 Instance (*class in nordlys.core.ml.instance*), 31  
 Instances (*class in nordlys.core.ml.instances*), 33  
 is\_name\_entity() (*in module nordlys.logic.el.el\_utils*), 53  
 is\_overlapping() (*nordlys.logic.el.greedy.Greedy method*), 54

**J**

JM (*nordlys.core.retrieval.scorer.ScorerLM attribute*), 46

**L**

LateFusionScorer (*class in nordlys.logic.fusion.late\_fusion\_scorer*), 62  
 len\_ratio() (*nordlys.logic.features.ftr\_mention.FtrMention method*), 60  
 link() (*nordlys.logic.el.cmns.Cmns method*), 53  
 link() (*nordlys.logic.el.ltr.LTR method*), 55

- link() (*nordlys.services.el.EL method*), 66
- LM\_MODELS (*nordlys.core.retrieval.retrieval.Retrieval attribute*), 44
- lm\_score() (*nordlys.logic.features.ftr\_entity\_similarity.FtrEntitySimilarity method*), 60
- load() (*nordlys.core.eval.trec\_qrels.TrecQrels method*), 28
- load\_associations() (*nordlys.logic.fusion.fusion\_scorer.FusionScorer method*), 62
- load\_config() (*nordlys.core.utils.file\_utils.FileUtils static method*), 52
- load\_entities() (*in module nordlys.logic.elr.field\_mapping*), 56
- load\_entities() (*in module nordlys.logic.er.field\_mapping*), 57
- load\_file() (*nordlys.core.eval.trec\_run.TrecRun method*), 29
- load\_folds() (*nordlys.core.ml.cross\_validation.CrossValidation method*), 31
- load\_kb\_snapshot() (*in module nordlys.logic.el.el\_utils*), 53
- load\_nordlys\_config() (*in module nordlys.config*), 69
- load\_queries() (*nordlys.logic.fusion.fusion\_scorer.FusionScorer attribute*), 56
- load\_results() (*nordlys.core.eval.trec\_eval.TrecEval method*), 27
- load\_yerd() (*nordlys.logic.el.ltr.LTR static method*), 55
- lookup\_en() (*nordlys.logic.entity.entity.Entity method*), 57
- lookup\_name\_dbpedia() (*nordlys.logic.entity.entity.Entity method*), 57
- lookup\_name\_facc() (*nordlys.logic.entity.entity.Entity method*), 57
- LTR (*class in nordlys.logic.el.ltr*), 54
- M**
- main() (*in module nordlys.core.data.dbpedia.indexer\_dbpedia\_types*), 23
- main() (*in module nordlys.core.data.facc.facc2mongo*), 24
- main() (*in module nordlys.core.data.word2vec.word2vec2mongo*), 24
- main() (*in module nordlys.core.eval.eval*), 25
- main() (*in module nordlys.core.eval.trec\_qrels*), 28
- main() (*in module nordlys.core.eval.trec\_run*), 29
- main() (*in module nordlys.core.ml.cross\_validation*), 31
- main() (*in module nordlys.core.ml.instance*), 33
- main() (*in module nordlys.core.ml.instances*), 34
- main() (*in module nordlys.core.ml.ml*), 37
- main() (*in module nordlys.core.retrieval.retrieval*), 45
- main() (*in module nordlys.core.retrieval.toy\_indexer*), 48
- main() (*in module nordlys.core.storage.mongo*), 51
- main() (*in module nordlys.core.storage.parser.nt\_parser*), 49
- main() (*in module nordlys.core.utils.file\_utils*), 52
- main() (*in module nordlys.logic.el.cmns*), 53
- main() (*in module nordlys.logic.elr.field\_mapping*), 56
- main() (*in module nordlys.logic.er.field\_mapping*), 57
- main() (*in module nordlys.logic.features.word2vec*), 61
- main() (*in module nordlys.logic.query.mention*), 63
- main() (*in module nordlys.services.ec*), 65
- main() (*in module nordlys.services.el*), 66
- main() (*in module nordlys.services.er*), 68
- main() (*in module nordlys.services.tti*), 69
- make\_plot() (*nordlys.core.eval.plot\_diff.QueryDiff method*), 26
- map() (*nordlys.logic.elr.field\_mapping.FieldMapping method*), 56
- map() (*nordlys.logic.er.field\_mapping.FieldMapping method*), 57
- MAPPING\_DEBUG (*nordlys.logic.elr.field\_mapping.FieldMapping attribute*), 56
- MAPPING\_DEBUG (*nordlys.logic.er.field\_mapping.FieldMapping attribute*), 57
- matches() (*nordlys.logic.features.ftr\_mention.FtrMention method*), 60
- mct() (*nordlys.logic.features.ftr\_entity\_mention.FtrEntityMention method*), 59
- Mention (*class in nordlys.logic.query.mention*), 63
- mention\_len() (*nordlys.logic.features.ftr\_mention.FtrMention method*), 61
- ML (*class in nordlys.core.ml.ml*), 36
- mlm\_score() (*nordlys.logic.features.ftr\_entity\_similarity.FtrEntitySimilarity method*), 60
- Mongo (*class in nordlys.core.storage.mongo*), 50
- multi\_termvector() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42
- N**
- name (*nordlys.core.data.dbpedia.indexer\_dbpedia\_types.IndexerDBpediaTypes attribute*), 23
- nordlys (*module*), 22
- nordlys.config (*module*), 69
- nordlys.core (*module*), 22
- nordlys.core.data (*module*), 22
- nordlys.core.data.dbpedia (*module*), 22
- nordlys.core.data.dbpedia.indexer\_dbpedia\_types (*module*), 23

nordlys.core.data.facc (*module*), 23  
 nordlys.core.data.facc.facc2mongo (*module*), 23  
 nordlys.core.data.word2vec (*module*), 24  
 nordlys.core.data.word2vec.word2vec2mongo (*module*), 24  
 nordlys.core.eval (*module*), 24  
 nordlys.core.eval.eval (*module*), 25  
 nordlys.core.eval.plot\_diff (*module*), 25  
 nordlys.core.eval.query\_diff (*module*), 26  
 nordlys.core.eval.trec\_eval (*module*), 26  
 nordlys.core.eval.trec\_qrels (*module*), 27  
 nordlys.core.eval.trec\_run (*module*), 28  
 nordlys.core.ml (*module*), 29  
 nordlys.core.ml.cross\_validation (*module*), 30  
 nordlys.core.ml.instance (*module*), 31  
 nordlys.core.ml.instances (*module*), 33  
 nordlys.core.ml.ml (*module*), 35  
 nordlys.core.retrieval (*module*), 37  
 nordlys.core.retrieval.elastic (*module*), 38  
 nordlys.core.retrieval.elastic\_cache (*module*), 42  
 nordlys.core.retrieval.indexer\_mongo (*module*), 43  
 nordlys.core.retrieval.retrieval (*module*), 43  
 nordlys.core.retrieval.retrieval\_results (*module*), 45  
 nordlys.core.retrieval.scorer (*module*), 45  
 nordlys.core.retrieval.toy\_indexer (*module*), 48  
 nordlys.core.storage (*module*), 48  
 nordlys.core.storage.mongo (*module*), 50  
 nordlys.core.storage.parser (*module*), 48  
 nordlys.core.storage.parser.nt\_parser (*module*), 48  
 nordlys.core.storage.parser.uri\_prefix (*module*), 49  
 nordlys.core.utils (*module*), 51  
 nordlys.core.utils.file\_utils (*module*), 51  
 nordlys.core.utils.logging\_utils (*module*), 52  
 nordlys.logic (*module*), 53  
 nordlys.logic.el (*module*), 53  
 nordlys.logic.el.cmns (*module*), 53  
 nordlys.logic.el.el\_utils (*module*), 53  
 nordlys.logic.el.greedy (*module*), 54  
 nordlys.logic.el.ltr (*module*), 54  
 nordlys.logic.elr (*module*), 55  
 nordlys.logic.elr.field\_mapping (*module*), 55  
 nordlys.logic.elr.top\_fields (*module*), 56  
 nordlys.logic.entity (*module*), 56  
 nordlys.logic.entity.entity (*module*), 56  
 nordlys.logic.er (*module*), 57  
 nordlys.logic.er.field\_mapping (*module*), 57  
 nordlys.logic.er.top\_fields (*module*), 58  
 nordlys.logic.features (*module*), 58  
 nordlys.logic.features.feature\_cache (*module*), 58  
 nordlys.logic.features.ftr\_entity (*module*), 59  
 nordlys.logic.features.ftr\_entity\_mention (*module*), 59  
 nordlys.logic.features.ftr\_entity\_similarity (*module*), 59  
 nordlys.logic.features.ftr\_mention (*module*), 60  
 nordlys.logic.features.word2vec (*module*), 61  
 nordlys.logic.fusion (*module*), 61  
 nordlys.logic.fusion.fusion\_scorer (*module*), 61  
 nordlys.logic.fusion.late\_fusion\_scorer (*module*), 62  
 nordlys.logic.query (*module*), 63  
 nordlys.logic.query.mention (*module*), 63  
 nordlys.logic.query.query (*module*), 63  
 nordlys.logic.tti (*module*), 64  
 nordlys.logic.tti.type\_centric (*module*), 64  
 nordlys.services (*module*), 64  
 nordlys.services.api (*module*), 64  
 nordlys.services.ec (*module*), 65  
 nordlys.services.el (*module*), 65  
 nordlys.services.er (*module*), 66  
 nordlys.services.tti (*module*), 68  
 normalize() (*nordlys.core.eval.trec\_run.TrecRun method*), 29  
 notanalyzed\_field() (*nordlys.core.retrieval.elastic.Elastic static method*), 40  
 notanalyzed\_searchable\_field() (*nordlys.core.retrieval.elastic.Elastic static method*), 40  
 NTParser (*class in nordlys.core.storage.parser.nt\_parser*), 49  
 num\_docs() (*nordlys.core.retrieval.elastic.Elastic method*), 40  
 num\_docs() (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42  
 num\_docs() (*nordlys.core.retrieval.retrieval\_results.RetrievalResults method*), 45  
 num\_fields() (*nordlys.core.retrieval.elastic.Elastic method*), 40

`num_fields()` (*nordlys.core.retrieval.elastic\_cache.ElasticCache* method), 42

`num_rel()` (*nordlys.core.eval.trec\_qrels.TrecQrels* method), 28

## O

`object()` (*nordlys.core.storage.parser.nt\_parser.Triple* method), 49

`object_prefixed()` (*nordlys.core.storage.parser.nt\_parser.Triple* method), 49

`OP_QUERY_DIFF` (*nordlys.core.eval.eval.Eval* attribute), 25

`open_file_by_type()` (*nordlys.core.utils.file\_utils.FileUtils* static method), 52

`OPERATIONS` (*nordlys.core.eval.eval.Eval* attribute), 25

`outlinks()` (*nordlys.logic.features.ftr\_entity.FtrEntity* method), 59

`output()` (*nordlys.core.ml.ml.ML* method), 36

## P

`parse_file()` (*nordlys.core.storage.parser.nt\_parser.NTParser* method), 49

`pos1()` (*nordlys.logic.features.ftr\_entity\_mention.FtrEntityMention* method), 59

`predicate()` (*nordlys.core.storage.parser.nt\_parser.Triple* method), 49

`predicate_prefixed()` (*nordlys.core.storage.parser.nt\_parser.Triple* method), 49

`print_doc()` (*nordlys.core.storage.mongo.Mongo* static method), 51

`print_stat()` (*nordlys.core.eval.trec\_qrels.TrecQrels* method), 28

`print_stat()` (*nordlys.core.eval.trec\_run.TrecRun* method), 29

`PrintHandler` (class in *nordlys.core.utils.logging\_utils*), 52

`properties` (*nordlys.core.ml.instance.Instance* attribute), 32

`prune_by_score()` (*nordlys.logic.el.greedy.Greedy* method), 54

`prune_containment_mentions()` (*nordlys.logic.el.greedy.Greedy* method), 54

## Q

`qid` (*nordlys.logic.query.query.Query* attribute), 63

`Query` (class in *nordlys.logic.query.query*), 63

`query` (*nordlys.core.retrieval.retrieval\_results.RetrievalResults* attribute), 45

`query` (*nordlys.logic.query.query.Query* attribute), 63

`QueryDiff` (class in *nordlys.core.eval.plot\_diff*), 25

## R

`rank_ens()` (*nordlys.logic.el.cmns.Cmns* method), 53

`rank_ens()` (*nordlys.logic.el.ltr.LTR* method), 55

`raw_query` (*nordlys.logic.query.query.Query* attribute), 63

`read_file_as_list()` (*nordlys.core.utils.file\_utils.FileUtils* static method), 52

`redirects()` (*nordlys.logic.features.ftr\_entity.FtrEntity* method), 59

`RequestHandler` (class in *nordlys.core.utils.logging\_utils*), 52

`Retrieval` (class in *nordlys.core.retrieval.retrieval*), 44

`retrieval()` (in module *nordlys.services.api*), 64

`RetrievalResults` (class in *nordlys.core.retrieval.retrieval\_results*), 45

`retrieve()` (*nordlys.core.retrieval.retrieval.Retrieval* method), 44

`retrieve()` (*nordlys.services.er.ER* method), 67

`Parse()` (*nordlys.core.eval.eval.Eval* method), 25

`run()` (*nordlys.core.ml.cross\_validation.CrossValidation* method), 31

`run()` (*nordlys.core.ml.ml.ML* method), 36

## S

`save_folds()` (*nordlys.core.ml.cross\_validation.CrossValidation* method), 31

`score_doc()` (*nordlys.core.retrieval.scorer.ScorerLM* method), 47

`score_doc()` (*nordlys.core.retrieval.scorer.ScorerMLM* method), 47

`score_doc()` (*nordlys.core.retrieval.scorer.ScorerPRMS* method), 48

`score_queries()` (*nordlys.logic.fusion.fusion\_scorer.FusionScorer* method), 62

`score_query()` (*nordlys.logic.fusion.fusion\_scorer.FusionScorer* method), 62

`score_query()` (*nordlys.logic.fusion.late\_fusion\_scorer.LateFusionScorer* method), 62

`Scorer` (class in *nordlys.core.retrieval.scorer*), 45

`SCORER_DEBUG` (*nordlys.core.retrieval.scorer.Scorer* attribute), 45

`ScorerLM` (class in *nordlys.core.retrieval.scorer*), 46

`ScorerMLM` (class in *nordlys.core.retrieval.scorer*), 47

`ScorerPRMS` (class in *nordlys.core.retrieval.scorer*), 47

`SCORES` (*nordlys.core.eval.plot\_diff.QueryDiff* attribute), 25

`search()` (*nordlys.core.retrieval.elastic.Elastic* method), 41

`search_complex()` (*nordlys.core.retrieval.elastic.Elastic* method), 41

set () (*nordlys.core.storage.mongo.Mongo method*), 51  
 set\_feature\_val () (*nordlys.logic.features.feature\_cache.FeatureCache method*), 58  
 SIMILARITY (*nordlys.core.retrieval.elastic.Elastic attribute*), 39  
 subject () (*nordlys.core.storage.parser.nt\_parser.Triple method*), 49  
 subject\_prefixed () (*nordlys.core.storage.parser.nt\_parser.Triple method*), 49  
 triple\_parsed () (*nordlys.core.storage.parser.nt\_parser.TripleHandler method*), 49  
 triple\_parsed () (*nordlys.core.storage.parser.nt\_parser.TripleHandler method*), 49  
 TripleHandler (*class in nordlys.core.storage.parser.nt\_parser*), 49  
 TripleHandlerPrinter (*class in nordlys.core.storage.parser.nt\_parser*), 49  
 TTI (*class in nordlys.services.tti*), 69  
 TypeCentric (*class in nordlys.logic.tti.type\_centric*), 64

## T

tcm () (*nordlys.logic.features.ftr\_entity\_mention.FtrEntityMention method*), 59  
 tem () (*nordlys.logic.features.ftr\_entity\_mention.FtrEntityMention method*), 59  
 term\_freq () (*nordlys.core.retrieval.elastic.Elastic method*), 41  
 term\_freq () (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 42  
 term\_freqs () (*nordlys.core.retrieval.elastic.Elastic method*), 41  
 term\_freqs () (*nordlys.core.retrieval.elastic\_cache.ElasticCache method*), 43  
 to\_elq\_eval () (*in module nordlys.logic.el.el\_utils*), 54  
 to\_json () (*nordlys.core.ml.instance.Instance method*), 32  
 to\_json () (*nordlys.core.ml.instances.Instances method*), 34  
 to\_libsvm () (*nordlys.core.ml.instance.Instance method*), 32  
 to\_libsvm () (*nordlys.core.ml.instances.Instances method*), 34  
 to\_str () (*nordlys.core.ml.instance.Instance method*), 32  
 to\_str () (*nordlys.core.ml.instances.Instances method*), 34  
 to\_treceval () (*nordlys.core.ml.instances.Instances method*), 34  
 TopFields (*class in nordlys.logic.elr.top\_fields*), 56  
 TopFields (*class in nordlys.logic.er.top\_fields*), 58  
 train () (*nordlys.logic.el.ltr.LTR static method*), 55  
 train\_model () (*nordlys.core.ml.ml.ML method*), 36  
 trec\_format () (*nordlys.core.retrieval.retrieval.Retrieval method*), 44  
 TrecEval (*class in nordlys.core.eval.trec\_eval*), 26  
 TrecQrels (*class in nordlys.core.eval.trec\_qrels*), 27  
 TrecRun (*class in nordlys.core.eval.trec\_run*), 28  
 Triple (*class in nordlys.core.storage.parser.nt\_parser*), 49  
 triple () (*nordlys.core.storage.parser.nt\_parser.Triple method*), 49

## U

unescape () (*nordlys.core.storage.mongo.Mongo static method*), 51  
 unescape\_doc () (*nordlys.core.storage.mongo.Mongo static method*), 51  
 update\_similarity () (*nordlys.core.retrieval.elastic.Elastic method*), 42  
 URIPrefix (*class in nordlys.core.storage.parser.uri\_prefix*), 49

## W

Word2Vec (*class in nordlys.logic.features.word2vec*), 61  
 Word2VecToMongo (*class in nordlys.core.data.word2vec.word2vec2mongo*), 24  
 write\_trec\_format () (*nordlys.core.retrieval.retrieval\_results.RetrievalResults method*), 45