

# **Dolphin Documentation**

*Release 0.0.1*

**Infodevelopers Pvt Ltd**

Jul 11, 2020



<b>Python Module Index</b>	<b>9</b>
<b>Index</b>	<b>11</b>



# dolphin

Dolphin is the **Dolphwin**'s Artificial Intelligence component – as intelligent as a Dolphin.

# Directory Structure

```
└─ dolphwin          # Main Directory
  └─ dolphin          # AI package for Dolphwin
    └─ archives
    └─ data            # data for pool parsing
      └─ languages.csv  # predefined pool for languages
      └─ nationality_data.csv # predefined pool for nationalities
      └─ soft_skills_resume.csv # predefined pool for soft skills
      └─ technical_skills.csv # predefined pool for technical skills
    └─ datasets        # data of resume, job description and other files
    └─ docs            # ? project document files
    └─ labs            # ? contains research program files
    └─ models          # machine learning models saved
      └─ NER_model.ser.gz # NER model saved from training with Stanford NLP
      └─ stanford_ner.jar # pretrained stanford model
      └─ segment_identifier.pkl # model trained for resume segmentation
      └─ new_model_svc.pkl # SVC model for document classification
      └─ tfidf1.pkl      # tfidf model used for document classification
    └─ segmentation    # module for resume segmentation
      └─ segment_resume.py # python module for resume segmentation
    └─ stanfordNER      # module for parsing using NER
      └─ formatdata.py   # module for formatting data
      └─ namedentityrecognition.py # module for getting named entities
      └─ standardizedata.py # module for standardizing data
    └─ datareader.py    # python module for reading textual content from multiple
      extension files
```

```
├── pipeline.py          # python module for mantaining every task in custom spacy
pipeline
├── evaluate_pipeline.py # python module for testing working of pipeline with a
test resume
├── settings.py         # python module for mantaining configurations
├── settings.json       # configuration in json format
├── tests               # all unit testing files
│   ├── test_datareader.py      # unit testing file for datareader module
│   ├── test_identify_resume.py # unit testing file for document categorization
module
│   ├── test_ner_parsing.py     # unit testing file for testing ner parsing module
│   └── test_segment_resume.py  # unit testing file for testing resume segmentation
module
├── setup.py            # setup file for pip
├── LICENSE             # License details
├── info.log            # Logs mantained using logging module
├── requirements.txt    # list of python modules required to run the system
└── README.md          # markdown documentation of the system
```

## # Pipelines

The pipelines used in this system are shown in the below diagram. ![pipeline](pipeline.png)

### ## Document Categorization pipeline

This is a custom spacy pipeline which is used for identifying the category of any textual documents. The input documents supports following file formats:

- .pdf
- .docx
- .doc
- .txt

The cateogory defined in the pipeline are as follows:

1. Resume
2. Job Decription
3. Others

This pipeline is developed training a Support Vector Classifier with training data containing 400 samples each of resume, job descriptions and other files.

> Trained model for document categorization —> new\_model\_svc.pkl

### ### Usage

```
``` from pipeline import nlp
doc = nlp(textual_content_of_file)
```

```
category = doc._category
```

```
'''
```

### ## Segmentation pipeline

This pipeline is used to segment any resume to following segments:

- Profile segment
- Objective segment
- SKills segment
- Academics segment
- Experience segment
- Language segment
- Projects segment
- Rewards segment
- References segment
- Links segment

If any document is categorized as resume, it is passed to the segmentation pipeline to break down into different resume segments.

> Trained model for resume segmentation -> segment\_identifier.pkl

### ### Usage

```
''' from pipeline import nlp
```

```
doc = nlp(textual_content_of_file)
```

```
profile_segment = doc._profile_segment objective_segment = doc._objective_segment skills_segment =  
doc._skills_segment academics_segment = doc._academics_segment experience_segment = doc._experi-  
ence_segment language_segment = doc._language_segment projects_segment = doc._projects_segment  
rewards_segment = doc._rewards_segment references_segment = doc._references_segmen links_seg-  
ment = doc._links_segment '''
```

### ## NER Parser pipeline

This pipeline is developed for the extraction of vital informations from any resume such as experience, education, personal information with name and address from any resume. For extractoin of such informa-tion from any resume, custom **Named Entity Recognition(NER)** model has been trained. This model has been trained using the **StanfordNLP** NER training interface. This model has been trained with the tokens of about 285 resumes.

The information extraction from resume using this pipeline are as follows:

- Name
- Address
- Education
- Experience

> Trained model for Named Entity Recognition —> NER\_model.ser.gz

### ### Usage

```
''' from pipeline import nlp
```

```
doc = nlp(textual_content_of_the_file) name = doc._.name address = doc._.address experience = doc._.experience education = doc._.education """
```

### ## Pattern Matching Pipeline

The information which follows a specific pattern or is contained with finite pool have been extracted using Spacy Matcher. The information extracted using this pipelines are as follows:

- technical skills
- soft skills
- nationality
- languages
- date
- phone number
- gender
- email

### ### Usage

```
""" from pipeline import nlp doc = nlp(textual_content_of_file)

nationality = doc._.nationality language = doc._.language email = doc._.email phone_number = doc._.phone technical_skills = doc._.technical_skills soft_skills = doc._.soft_skills date = doc._.date gender = doc._.gender """
```

### # Modules Used

- PyMuPDF
- texteract
- datefinder
- docx2xt
- gensim
- nltk
- spacy
- pytesseract
- sklearn

> Other requirement —> Downloading en\_core\_web\_sm modle from spacy —> python -m spacy download en\_core\_web\_sm

```
dolphin.datareader.clean_text ( text, dolower )
```

**Accepts the plain text and makes use of regex for cleaning the noise** :param: text :type:str :return:- cleaned text :type str

```
dolphin.datareader.doc_to_text ( filepath, dolower )
```

**Takes the doc file from the file path param and returns the cleaned the text from the file.** :param filepath: path/directory of the doc file in the system :return: Returns the cleaned text from the file

```
dolphin.datareader.docx_to_text ( file_path, dolower )
```

**Takes docx files and extracts plain text from the docx files** :param file\_path :type str :return:text :type str



`dolphin.datareader.img_to_text ( filepath, dolower )`

**Takes the image file from the file path param and returns the cleaned the text from the image file.**  
 :param filepath: path/directory of the image file in the system :return: Returns the cleaned text from the image file

`dolphin.datareader.pdf_to_text ( file_path, dolower )`

**Takes filepath and extracts the plain text from pdf for training the word to vec model** :param file\_path :type str :return:text :type str

`dolphin.datareader.prepare_text ( file, dolower )`

**Takes the resume or any other doc; checks the extension of doc and then uses suitable methods to extract and clean the text** :param: file :type str :return: cleaned tokenized sentences :type list

`dolphin.datareader.prepare_text_from_string ( text, dolower )`

**takes string of text and then cleans the noise** :param: text :type str :return: cleaned\_text :type str

`dolphin.datareader.txt_to_text ( file_path, dolower )`

**Extracts plain text from txt files** :param file\_path :type str :return:text :type str

**class** `dolphin.document_categorization.doc_classifier_test.DocClassifier` ( `tfidfmodelpath, classifiermodelpath, nlp` )

**This class is developed for the preprocessing of textual contents in any document and use such preprocessed content for the prediction of category of the document using tfidf vectorizer developed and SVC classifier model developed**

**Param** tfidfmodelpath: The path of the tfidf vectorizer developed from the training data

**Type** str

**Param** classifiermodelpath: The path of the classifier model developed by training SVC classifier

**classify** ( `file_content` )

**Function to apply preprocessed text to the document categorization model** :param: file\_content: textual file content of the document :type: str :return: category of the document :rtype: str

**preprocess\_words** ( `texts, stop_words, bigram_mod, trigram_mod, allowed_postags=['NOUN', 'ADJ', 'VERB', 'ADV']` )

**Remove Stopwords, Form Bigrams, Trigrams and Lemmatization** :param: texts: list of tokens in the document type: list :param: stop\_words: list of stopword from nltk corpus :type: list :param: bigram\_mod: bigrams from document evaluated using gensim Phraser :type: list :param: trigram\_mod: trigrams from document evaluated using gensim Phraser :type: list :param: allowedpos\_tags: list of allowed POS tags in the document :type: list :return: preprocessed list of tokens from document :rtype: list

**sent\_to\_words** ( `file_content` )

**remove emails and preprocess the sentences in the document using gensim simple\_preprocess function** :param: file\_content: the textual content of the document :type: str :return: list of preprocessed sentences :rtype: list

**class** `dolphin.segmentation.segmentresume.ResumeSegmentCreator`

**This class is developed for the segmentation of resumes into several segments. The segments are personal info, objective, skills, experiences, projects, academics, rewards, languages, references and links**

**format\_segment** ( `text` )

**This function formats the resume segments** :param: text : textual content of the resume : type:

str :return: resume segments extracted :rtype: dict

**sliced\_resume\_text ( )**

This function returns the sliced resume text return: sliced text rtype: dict

**unique\_index\_headings ( )**

This function extracts the unique index headings from the resume content return: unique index headings rtype: list

dolphin.stanfordNER.formatdata.**formatPersonalinfo ( personal\_info )**

This functions takes all the unstructured personal information and returns the formatted and structured personal information. :param personal\_info: (List) of all the fields required as personal info. :return: List(Returns the formatted personal info will all the required fields.)

**class dolphin.stanfordNER.namedentityrecognition.StanfordNER**

This class applies the trained custom NER model to return the academics, experiences and personal info from any resume

**static education\_parser ( tagged\_tuples, text )**

**\*\* Function to gather all the education predicted tokens and format the tokens\*\*** param: tagged\_tuples: tagged\_tuples from NER model type: list of tuples param: text: textual content to be applied for education parsing type: str

**static experience\_parser ( tagged\_tuples, text )**

**Function to get out the experience related tokens from NER mmodel and format the parsed information** param: tagged\_tuples: tagged\_tuples from NER model type: list of tuples param: text: textual content to be applied for experience parsing type: str return: formatted\_data rtype: dict

**static ner\_parser ( model, text, mode )**

**NER model that chooses the particular model and parser** param: model: Stanford NLP model param: text: textual content of the resume type: str param: mode: to select the respective parser type: str

**static personal\_info\_parser ( tagged\_tuples, text )**

**Function to be applied for the parsing of personal information from resumes** param: tagged\_tuples: tagged\_tuples from NER model type: list of tuples param: text: textual content to be applied for experience parsing type: str

dolphin.stanfordNER.namedentityrecognition.**gather\_ner ( tagged\_tuples, ner='DESIG' )**

**This method gathers the predicted tokens related to academics, experiences and personal info from NER** :param:tagged\_tuples :list of tuples :return : list of entities :rtype:list

dolphin.stanfordNER.standarizedata.**capitalizeinput ( inputparameter )**

This will take the raw input text and tokenize the string and capitalize all the token and returns the formatted capitalized text. :param inputparameter :type str :return: string(text will all the tokens capitalized)

dolphin.stanfordNER.standarizedata.**comparedates ( datelist )**

this function takes the list of the date parse them and compares them and provide larger date and smaller date :param datelist:type list of str :return: sorted dates

**class dolphin.pipeline.NlpPipeline**

This class is developed to add the custom spacy pipelines for document categorization, resume

segmentation, and resume parsing.

**add\_category\_component ( )**

Function to add document categorization component to the spacy custom pipeline

**add\_ner\_parsing ( )**

Function for adding custom NER pipeline in nlp model

**add\_pattern\_matching ( )**

Function for adding pattern matching component in pipeline

**add\_segmentation\_component ( )**

Function to add resume segmentation component to the spacy custom pipeline

**category\_component ( doc )**

custom function to add resume identification component to the pipeline params: spacy Doc of textual data return : Doc object with category embedded

**ner\_parsing\_component ( doc )**

Function for preparing custom ner parsing component params: doc of content from file return: doc containing parsed information from keyword

**pattern\_matching\_component ( doc )**

Function for preparing pattern matching component in spacy params: doc from document return : doc containing matched keywords

**process\_text ( content )**

Function to process as spacy doc

**segmentation\_component ( doc )**

custom function to add resume identification component to the pipeline params: spacy Doc of textual data return : Doc object with category embedded

- [Index](#)
- [Module Index](#)
- [Search Page](#)



## d

dolphin

    dolphin.datareader, ??

    dolphin.document\_categorization.doc\_classifier\_test,  
        ??

    dolphin.pipeline, ??

    dolphin.segmentation.segmentresume,  
        ??

    dolphin.stanfordNER.formatdata, ??

    dolphin.stanfordNER.namedentityrecognition,  
        6

    dolphin.stanfordNER.standarizedata,  
        6



## A

add\_category\_component() (dolphin.pipeline.NlpPipeline method), 7  
 add\_ner\_parsing() (dolphin.pipeline.NlpPipeline method), 7  
 add\_pattern\_matching() (dolphin.pipeline.NlpPipeline method), 7  
 add\_segmentation\_component() (dolphin.pipeline.NlpPipeline method), 7

## C

capitalizeinput() (in module dolphin.stanfordNER.standarizedata), 6  
 category\_component() (dolphin.pipeline.NlpPipeline method), 7  
 classify() (dolphin.document\_categorization.doc\_classifier\_test.DocClassifier method), 5  
 clean\_text() (in module dolphin.datareader), 4  
 comparedates() (in module dolphin.stanfordNER.standarizedata), 6

## D

doc\_to\_text() (in module dolphin.datareader), 4  
 DocClassifier (class in dolphin.document\_categorization.doc\_classifier\_test), 5  
 docx\_to\_text() (in module dolphin.datareader), 4  
 dolphin.datareader (module), 4  
 dolphin.document\_categorization.doc\_classifier\_test (module), 5  
 dolphin.pipeline (module), 6  
 dolphin.segmentation.segmentresume (module), 5  
 dolphin.stanfordNER.formatdata (module), 6  
 dolphin.stanfordNER.namedentityrecognition (module), 6  
 dolphin.stanfordNER.standarizedata (module), 6

## E

education\_parser() (dolphin.stanfordNER.namedentityrecognition.StanfordNER static method), 6  
 experience\_parser() (dolphin.stanfordNER.namedentityrecognition.StanfordNER static method), 6

## F

format\_segment() (dolphin.segmentation.segmentresume.ResumeSegmentCreator method), 5  
 formatPersonalinfo() (in module dolphin.stanfordNER.formatdata), 6

## G

gather\_ner() (in module dolphin.stanfordNER.namedentityrecognition), 6

## I

img\_to\_text() (in module dolphin.datareader), 5

## N

ner\_parser() (dolphin.stanfordNER.namedentityrecognition.StanfordNER static method), 6  
 ner\_parsing\_component() (dolphin.pipeline.NlpPipeline method), 7  
 NlpPipeline (class in dolphin.pipeline), 6

## P

pattern\_matching\_component() (dolphin.pipeline.NlpPipeline method), 7  
 pdf\_to\_text() (in module dolphin.datareader), 5  
 personal\_info\_parser() (dolphin.stanfordNER.-

namedentityrecognition.StanfordNER  
static method), 6  
prepare\_text() (in module dolphin.datareader), 5  
prepare\_text\_from\_string() (in module dolphin.-  
datareader), 5  
preprocess\_words() (dolphin.document\_catego-  
rization.doc\_classifier\_test.DocClassifier  
method), 5  
process\_text() (dolphin.pipeline.NlpPipeline  
method), 7

## R

ResumeSegmentCreator (class in dolphin.seg-  
mentation.segmentresume), 5

## S

segmentation\_component() (dolphin.pipe-  
line.NlpPipeline method), 7  
sent\_to\_words() (dolphin.document\_categoriza-  
tion.doc\_classifier\_test.DocClassifier  
method), 5  
sliced\_resume\_text() (dolphin.segmentation.seg-  
mentresume.ResumeSegmentCreator  
method), 6  
StanfordNER (class in dolphin.stanfordNER.-  
namedentityrecognition), 6

## T

txt\_to\_text() (in module dolphin.datareader), 5

## U

unique\_index\_headings() (dolphin.segmenta-  
tion.segmentresume.ResumeSegmentCreator  
method), 6