

COMP4650 / COMP6490

Document Analysis 2018

Assignment Project Exam Help

Info

ction

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Gabriela Ferraro



Australian
National
University

Overview of IE lectures

- Introduction to Information Extraction (IE)

Overview Assignment Project Exam Help

Relation Extr

<https://eduassistpro.github.io/>

Named Entit

Add WeChat edu_assist_pro

- Sequence labeling method 2

- Automatic Summarization

* Acknowledgement: Some of the content originates from the Stanford NLP course at Coursera.org

Books

Speech and Language Processing

Jurafsky and Martin

2014. Pearson.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Na *age Processing*

Jacob Eisenstein

2018. MIT pres.

<https://github.com/jacobeisenstein/gt-nlp-class>

Introduction to IE

What is IE?

Automatically extract structured information from unstructured and/or semi-structured data.

Assignment Project Exam Help

Who did what <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Main goals:

- Helps natural language understanding
- Organize information for humans
- Organize information in a formal and precise form that allows further analysis and/or inferences made by computer algorithms

IE Applications

Scan documents and populate:

Templates

Ontologies

Data Bases

Knowledge Bases <https://eduassistpro.github.io/>

[Add WeChat edu_assist_pro](#)

Text understanding (e.g.: named entity recognition, relation extraction)

Automatic summarization

Question answering

...

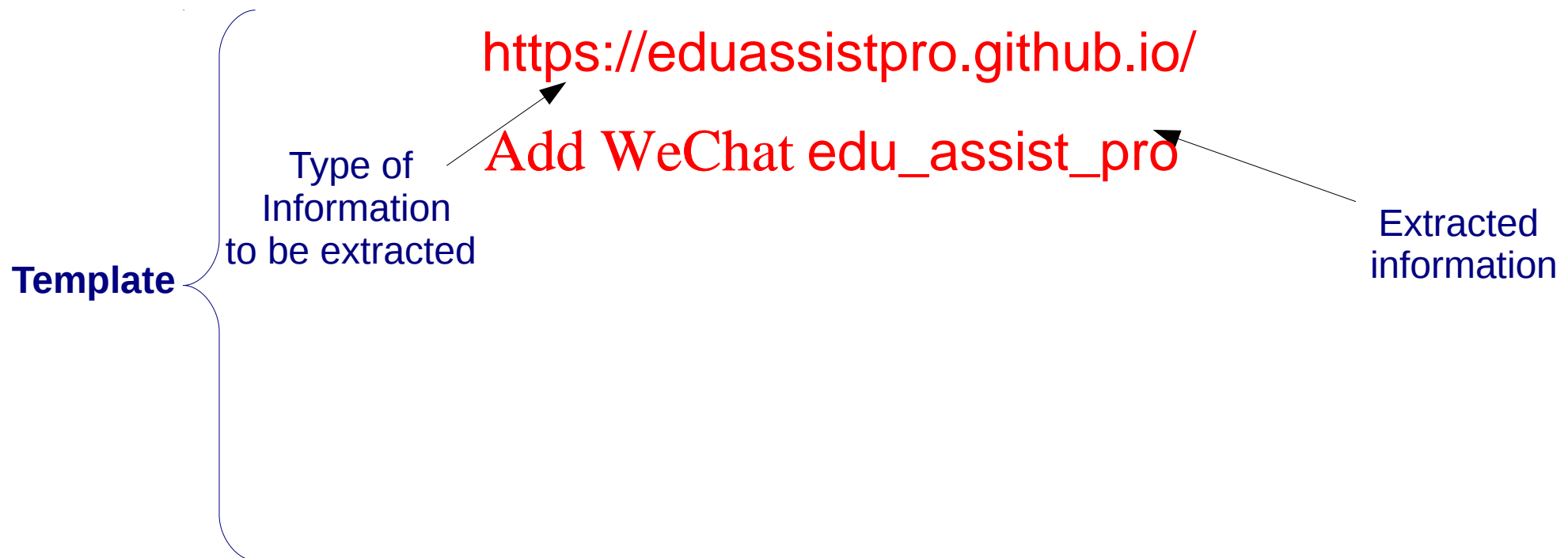
etc.

IE Template based example

2008 January 17

British Airways Flight 38, a **Boeing 777- 200ER**, lands short of the runway at **London Heathrow Airport** in the United Kingdom. **Nine** of the 152 people on board are treated for minor injuries, but there are **no fatalities**; this is the first loss of a Boeing 777.

Assignment Project Exam Help



Extract information about aircraft accidents from news

Templates types

Slots in a template are usually filled by a substring of a document

Assignment Project Exam Help

→ Some slots may

fillers

hysic Job type: n

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

→ Some slots may allow **multiple fillers**

Programming language: Java, C++, Python, etc.

IE applications

- **Relation Extraction**

Paris **is the capital** of France.

France's **capital** is Paris.

Paris <is-a-capital-of> France

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- **Name Entity Recognition**

Paris is the capital of

<LocationEntity> <...> <LocationEntity>

Add WeChat edu_assist_pro

- **Combined**

<LocationEntity> <is-a-capital-of> <LocationEntity>

IE methods

Hand written patterns

Supervised machine learning

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Semi-supervised and unsupervised learning

Add WeChat edu_assist_pro

Relation Extraction

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

How relations are express in natural language?

- Relations are instantiated by predicates
- Predicates have arguments
- Verbs are the most productive predicate form

<https://eduassistpro.github.io/>

`predicate (arg1, arg2)`

Mery likes cake.

likes (Mery, cake)

*Mery **rent** a boat for 2 weeks for 300 dollars.*

```
rent (Mery, boat, 2 weeks, 300 dollars)
```

Why Relation Extraction?

Create new structured knowledge, e.g., facts

- Augment current knowledge bases

Adding words, facts to
FreeBase or

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Support question answering

Which actor starred in the film BATMAN 3?

acted-in(?x, BATMAN)

is-a(?y, actor)

But which relations should we extract? And how?

Which relations to extract?

- A pre-defined set of relations
- All relations (e.g., all verbs and their arguments)
- Ontological relations

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Example of a pre-defined set of relations

17 relations from SemE 2008 “Relation Extraction Task

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Example of extracting ALL relations

Use syntactic dependency trees to extract predicates and their arguments

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Ontological relations

Examples from the WordNet Thesaurus <http://wordnetweb.princeton.edu/perl/webwn>

Hypernym (is-a): subsumption between classes

- Giraffe **IS--A** ruminant **IS--A** ungulate **IS--A** mammal **IS--A** vertebrate **IS--A** animal...

Assignment Project Exam Help

Hyponym relation or ^{an individual and} class

<https://eduassistpro.github.io/>

- Dog → Terrier → Bull Terrier → St. Bernard → ...
- San Francisco **instance-of** city

Add WeChat edu_assist_pro

Synonym relation

- **Car** Sense 1 => auto, automobile, motorcar, machine
- **Man** Sense 1 => adult men
- **Man** Sense 2 => homo, human being, human

Relation extraction projects

Resource Description Framework (RDF) triples

Golden Gate Park **location** San Francisco

Dbpedia: +1 billion eduassistpro.github.io/

dbpedia:Golden_Gate <https://eduassistpro.github.io/>

dbpedia:San_Francisco [Add WeChat edu_assist_pro](https://eduassistpro.github.io/)

Freebase relations: well-known people, places, and things

<https://www.freebase.com/>

Total RDF triples: 2.1M

How to build relation extractors

Hand written patterns

Supervised machine learning

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Semi-supervised and unsupervised learning

Add WeChat edu_assist_pro

Hand written rules: Hearst's Patterns for extracting IS-A relations

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Marti A. Hearst. 1992. Automatic acquisition of hyponyms from large text corpora. In Proceedings of the 14th conference on Computational linguistics - Volume 2 (COLING '92), Vol. 2. Association for Computational Linguistics, Stroudsburg, PA, USA, 539-545.

Extracting Richer Relations Using Rules

Intuition: relations often holds between specific entities

- **located**(ORGANIZATION, LOCATION)
- **founded**(PERSON, <https://eduassistpro.github.io/>)
- **cures**(DRUG, DISEASE)

Start with Named Entity tags to help relation extraction

Which relations hold between 2 entities?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Which relations hold between 2 entities?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Summary: Hand written patterns for Relation Extraction

- **Plus**

- Human patterns tend to be high-precision
- Can be tailored to specific domains

- **Minus**

- Human patterns are often noisy
- A lot of work to think of all possible patterns
- Don't want to have to do this for every relation
- We would like better accuracy

Supervised machine learning for Relation Extraction

Training

- Choose the set of relations you want to extract
- Find and label data* = training set creation
- Extract relevant training set
- Train a classifier on the training set

Testing

- Tuned the classifier parameters on the dev. set
- Test the classifier on the test set

* Available RE datasets: SemEval7; BioNLP, etc.

Supervised relation extraction between entities

- Find all pairs of named entities (person, location, organization)
 - Decide if 2 entities are related
 - If yes, classify the types (is-a, instance-of)
- You can use any classifier you like
 - MaxEnt, Naive Bayes, CRF, SVM, CNN, etc.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

RESIDE: Improving Distantly-Supervised Neural Relation Extraction using Side Information

(Vashishth et al., 2018)

Summary: Supervised machine learning for Relation Extraction

- **Plus**

- Can get high accuracy with enough hand-labeled training data, if test data is similar enough to training data

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- **Minus**

Add WeChat edu_assist_pro

- Labeling a large training set is expensive
- Supervised models are brittle, don't generalize well to different genres

Semi supervised Relation Extraction

No training set? Maybe you have:

- A few **seed** tuples

- A few high-pre

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Can you use those **seeds** to do something useful?

- Bootstrapping: use the seeds to directly learn to populate a relation

Relation Bootstrapping (Hearst, 1992)

- Gather a set of **seed pairs** that have relation ***R***
- Iterate:
 - Find sentence **Assignment Project Exam Help**
<https://eduassistpro.github.io/>
 - Look at the context between **Add WeChat edu_assist_pro** and the pair and generalize the context to create patterns
 - Use the patterns for *grep* for more pairs

Bootstrapping

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Unsupervised relation extraction

- Extract relations from with no training data, thus no pre-defined list of relations

Assignment Project Exam Help

- Single-past: extra NPs

<https://eduassistpro.github.io/>

- Assessor ranks relations based on frequency

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Evaluation of unsupervised relation extraction

- Since it extracts totally new relations...

there is no gold set of correct relations

- cannot compute precision (don't know which ones are correct)
- cannot compute recall (don't know which ones were missed)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Instead, we can approximate **precision**

draw a random sample of relation from output, check precision manually

Name Entity Recognition

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Name Entity Recognition

Named Entity Recognition (NER)

Find and classify names in texts,
e.g.: person, location, organization, number, currency, etc.

Assignment Project Exam Help

Designated **S/2004** <https://eduassistpro.github.io/> own moon
to circle the giant planet. Add WeChat edu_assist_pro

It also appears to be the smallest moon in the
Neptunian system, measuring just **20 km (12 miles)**
across, completing one revolution around **Neptune**
every **23 hours**.

US astronomer **Mark Showalter** spotted the tiny dot
while studying segments of rings around Neptune.

proper name
quantity
location
person
Time
Other

NER Applications

- Machine Translation
- Question Answering **Assignment Project Exam Help**
- Automatic Summarization **<https://eduassistpro.github.io/>**
- Relation Extraction **Add WeChat edu_assist_pro**

NER as learning

- **Training**

- Collect a set of representative training documents
- Label each token for its entity class or other
- Design feature classes to the text and
<https://eduassistpro.github.io/>
- Train a sequence classifier to labels from the data
[Add WeChat edu_assist_pro](#)

- **Testing**

- Receive a set of testing documents
- Run sequence model inference to label each token
- Appropriately output the recognized entities

NER Task: the training data

US

astronomer

Mark

Showalter

spotted

that

□

□

LOC

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro,

O

O

□

□

Standard evaluation is per am Help

Entity not per token

precision, recall and F-measure

NER Task: example features

Numbers

- twoDigitNum (90) = Two-digit year
- fourDigitNum (1990) = Four-digit year
- containsDigitAndDash (09-96)
- containsDigitAndSlash (11/9/8)
- containsDigitAndComma (23,000.00) = Monetary amount
- containsDigitAndPeriod (1.00) = Monetary amount, percentage

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

NER Task: example features

- Person

- capPeriod (M.) = Person name initial
 - initCap (Sally d
 - lowerCase (c word
- <https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

- Organization

- allCaps (IBM) = Organization

* **Gazzeters** (list with persons, organizations, abbreviations, etc.)

NER challenges

Ambiguity problems:

- **Paris** (city vs. person)
 - **May** (person vs. month)
 - **2013** (date vs. quantity)
 - **Ferrari** (person vs. car)
- Assignment Project Exam Help
<https://eduassistpro.github.io/>

Multi-language NER:

- Language independent features (position, suffix, prefix, digits, POS-tags)
 - Lack of capitalization (Chinese, Indian lang., etc.)
 - Too much capitalization (German)
 - Free word order languages (Hungarian, Russian, etc.)
 - Languages with rich morphology (Czech, Spanish, etc.)
- Add WeChat edu_assist_pro

Evaluation in IE

How much relevant information has been extracted

Precision = $\frac{\text{\# of correct answers given by the system}}{\text{total \# of possible correct answers in the text}}$

Assignment Project Exam Help
How much of the extracted information is correct

Recall = $\frac{\text{\# of correct answers}}{\text{\# of answers given}}$
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

How good is the system in ignoring spurious information

Fall out = $\frac{\text{\# of incorrect answers given by the system}}{\text{\# of spurious facts in the text}}$

Combination of Precision and Recall

F-Measure = $2 * (\text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall}))$

IE take away

IE deals with processing human language texts by means of natural language processing techniques

- **Rule based methods**

- Use lexical patterns, e.g.: *X was born in Y.*
- Use syntactic patterns

- **Supervised method**

- Sequential labeling algorithms as HMM, CRF
- Required training data

- **Semi-supervised and unsupervised methods**

- Semi: required seed examples, e.g. lexical patterns
- Unsupervised: require unlabeled data
- Evaluation is not straightforward

Conclusion

- In the future, IE from cross-website pages will become more important as we move towards the Semantic
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- IE new challenges are: domain independent solutions, data integration and multilingualism
 - Lots need to be done!

Resources/Tools

KnowItAll

<https://github.com/knowitall>

Assignment Project Exam Help

**Stanford Named Entity
Pereira, 2001)**

<https://eduassistpro.github.io/>, McCallum, and

Add WeChat edu_assist_pro

<http://nlp.stanford.edu/software/CR>

OpenIE

<https://nlp.stanford.edu/software/openie.html>