An organizational approach for discipline prediction in research projects

Hoang-Son Pham

ECOOM UHasselt

EuroCRIS membership meeting, May 31, 2023, Brussels





- Proposed approach
- 3 Preliminary results







イロト イヨト イヨト





Hoang-Son Pham

Discipline prediction

)23

イロト イヨト イヨト イヨト

3/13

Classifying research documents is an essential task





Classifying research documents is an essential task

Automatically classifying disciplines is challenging





Classifying research documents is an essential task

Automatically classifying disciplines is challenging

Current approaches: citation-based, machine learning algorithms



Project metadata

- title
- keywords
- abstract

- researchers
- organisations
- disciplines



イロト イヨト イヨト

Project metadata

- title
- keywords
- abstract

- researchers
- organisations
- disciplines

Machine learning approach for discipline prediction



Predict disciplines of project based on disciplines of researchers

Researcher disciplines

- profile
- organizations
- projects
- co-authors on projects
- publications
- co-authors on publications



Each researcher is represented by a matrix $(N \times 6)$

	profile	organization	projects	co-authors on projects	publications	co-authors on publications
discipline 1	0.1	0.2	0.3	0.1	0.1	0.2
discipline 2	0.2	0.3	0	0	0.2	0.1
discipline 3	0.3	0.4	0.2	0.3	0	0
discipline N						

where a cell (i,j) is frequency of discipline i in resource j



► UHASSELT

Each researche	r is represented	by a matrix	$(N \times 6)$
----------------	------------------	-------------	----------------

	profile	organization	projects	co-authors on projects	publications	co-authors on publications
discipline 1	0.1	0.2	0.3	0.1	0.1	0.2
discipline 2	0.2	0.3	0	0	0.2	0.1
discipline 3	0.3	0.4	0.2	0.3	0	0
discipline N						

where a cell (i,j) is frequency of discipline i in resource j

A project is represented by a matrix (N \times 6), i.e., sum of researchers matrices



(日) (四) (日) (日) (日)

IIHASSEI

Each researche	er is	represented	by	а	matrix	(N	x 6)
----------------	-------	-------------	----	---	--------	----	-----	---

	profile	organization	projects	co-authors on projects	publications	co-authors on publications
discipline 1	0.1	0.2	0.3	0.1	0.1	0.2
discipline 2	0.2	0.3	0	0	0.2	0.1
discipline 3	0.3	0.4	0.2	0.3	0	0
discipline N						

where a cell (i,j) is frequency of discipline i in resource j

A project is represented by a matrix (N \times 6), i.e., sum of researchers matrices

Project data for training models is a 3-dimension matrix (M × N × 6)

Deep learning models

• Recurrent neural networks

Long Short-Term Memory Networks (LSTM)

Bidirectional Long Short-Term Memory Networks (BiLSTM)



Deep learning models

• Recurrent neural networks

Long Short-Term Memory Networks (LSTM)

Bidirectional Long Short-Term Memory Networks (BiLSTM)

• Convolutional Neural Networks (CNN)



Experimental setup

Dataset: 3954 projects, 42 disciplines Feature matrix: (3954, 42, 6) Training and testing: 80%, 20%





Hoang-Son Pham

UHASSEL'

Experimental setup

• Model hyperparameters

- numbers of hidden units = 128
- epochs = 100
- optimizer='adam'



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Experimental setup

Model hyperparameters

- numbers of hidden units = 128
- epochs = 100
- optimizer='adam'

Metrics

- accuracy
- precision
- recall
- f1-score
- jaccard index





Performance of models





23

Image: A match a ma

Performance of the models on the dataset, after excluding projects containing low-frequency disciplines



PERTISECENTRUM OLO MONITORING

HASSEI

We proposed a new approach to predict disciplines related to a project, based on the disciplines of researchers

Performance of deep learning models was relatively high

Further data pre-processing could improve the models' performance



► UHASSELT

We proposed a new approach to predict disciplines related to a project, based on the disciplines of researchers

Performance of deep learning models was relatively high

Further data pre-processing could improve the models' performance

Future work

- Further evaluate models on larger datasets
- Experiment with different hyperparameters
- Consider preprocessing the data differently, such as using different normalization techniques or feature engineering.



Thank you for listening



► UHASSELT

Hoang-Son Pham

Discipline prediction

)23

イロト イヨト イヨト