

Revealing Master's theses structures using Machine Learning Methods

Susman M., Hint H., Šinkūnienė J., Leijen D. A. J.

BURITE

Iceland
Liechtenstein
Norway grants



The Bwrite Project

Academic Writing in the Baltic States: Rhetorical Structures through culture(s) and languages

Goal

- **Measure** and map the writing traditions of Estonian, Latvian and Lithuanian
- **Develop a research method** to determine which features of a text are related to genre, discipline, culture and experience
- **Provide empirical results** that allow writers and instructors of writing to better apply those text features for teaching and writing

Website: <https://www.bwrite.ut.ee/>

Introduction

Academic texts are expected to follow an organizational structure

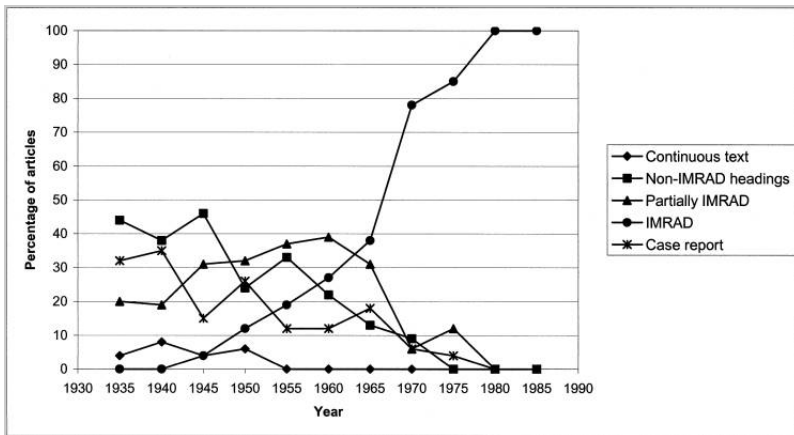
IMRaD structure as a standard in some communities (English, medicine, STEM, . . .)

Introduction, **M**ethod, **R**esults, and **D**iscussion

Research article genre, often imposed by the journal

Lin & Evans 2012

“IMRD is not an especially prevalent pattern in contemporary RA writing, so strict adherence to such a structure when conducting move-based or linguistic analyses is likely to result in incomplete or unrepresentative findings.”



Sollaci, L. B., & Pereira, M. G. (2004). The introduction, methods, results, and discussion (IMRAD) structure: a fifty-year survey. *Journal of the Medical Library Association : JMLA*, 92(3), 364–367.

Questions

- What do we observe in other, longer academic text genres, e.g. in master's theses? Can we detect IMRaD?
- What do we observe in languages other than English?
- What other structures the algorithm can detect?
- (What structures or patterns do non-expert writers follow?)

Manual annotations

Our database: web-scraped academic texts (PDFs), all available genres, disciplines, journals

Random selection of 467 theses in Lithuanian and Estonian

Based on TOC's, manually annotated on a 4-way continuum:

- Not IMRaD (NI),
- Rather Not IMRaD (RNI),
- Rather IMRaD (RI),
- IMRaD (I)

Overview of the method

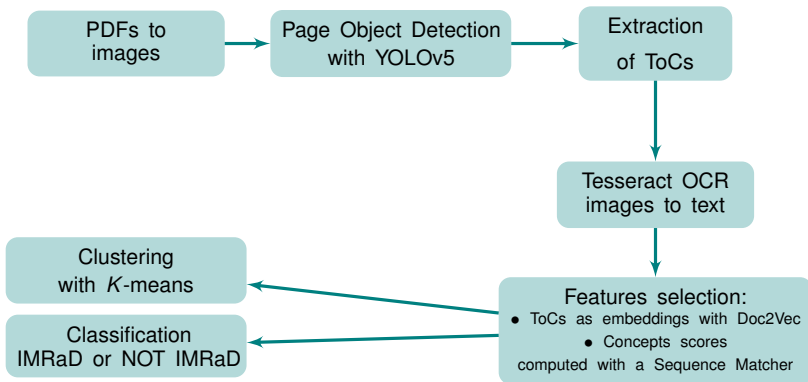
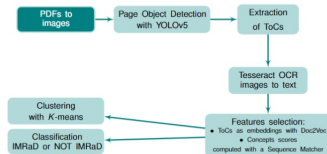


Figure: Revealing the structure of Masters theses: the method's steps

PDFs as images, why?

- Avoid loss of visual information
- Ease of extraction of relevant features



Esilehekülge

1 TARTU ÜLDKOOL FIILOSOOFIALEADUSKOND GERMAANI-ROMAANI FIILOLOGIA OSAKOND Meesko, Jekke HAIGE UURIMISE SAKSA-EESTI SÕNASTIK Registreeritud juhendaja: Anne Aroid, PhD TARTU SISUKORD Sissejuhatus .

2 Sõnastiku sisse ja allikate valik .

3 Eesti oskussõnastik .

4 Patsiendi uurimine .

5 Anamnees, läbivaatus ja instrumentaalne uuring .

6 Organismi talitlus ja elundite toimimine .

7 Varasem tervislik seisund ja pöördumised .

8 Sotsiaalsed olud ja eluviisid .

9 Kliiniline uuring .

10 Staatuse .

11 Uurimisjärjekord .

12 Haigusloog .

13 Terminoloogia ja termin .

14 Meditsiini oskussõnad .

15 Eesti meditsiini terminoloogia .

16 Sõnastiku koostamise tekkimise probleemid .

17 Sõnastiku koostamise heitumise uurimise saksakeelsete lühendid Eesti-saksa registreeritud kasutatud kirjandus Resümee Lisas SISSEJUHATUS Viimasele ajal on järjest avardunud võimalused teiste riikide, praktikate või tööde nimel, seda välisloog seaduse Euroopa liidu liikumisele.

18 14. Haiguste uurimise suundade vahetuseprobleemide rühmas, sissejuhatuse ja teostuste abil või enda ettevalmistel välismaale oma erialaseid teadmisi täiendada.

19 Tänu liikumisele on Eesti muutunud tulumaks ja atraktiivsemaks ning peetud on suund ke vastupidine - kasvades on siin õppivate, töötavate ja pühast veevate välismaalaste osakaal.

20 Arstidel ja meditsiiniüliõpilastel, kes on saksakeelses riigis tööl või õppimas, tuleb tihti kokku püüda haige uurimise ja küsitlusega.

21 Samuti võivad reisil viibivad isikud, kes ei pruugi osata antud maal räägitavat keelt, haigestuda või õnnetusse sattuda ning on sel juhul sunnitud arsti poole pöörduma.

22 18. Et saada teada, mis patsient on talin vastavõttel, panna diagnoos ja määrata ravi, on vajalik põhjalik küsitlus ning praktiline läbivaatus.

23 Anamnees ja läbivaatus on traditsioonilised võtted, mida peab valdama iga meditsiiniüliõpilane, resident ja arst.

24 Eriti olulisele vastuvõtul küsitatakse ja uuritakse patsienti üksikasjalikult ning sageli võtavad saksia ja eestlased arsti eluajaks anamneesi- ja läbivaatusküsimustik, milles nad täidavad saadud vastuste põhjal lühid.

25 Ühest poolt aitab see säästa arsti aega patsiendi uurimisel ja on abiks, et ta anamneesi võtmisel midagi ära ei unustaks, samuti võib vastuste käigus selguda haiguste raskus, millele midagi ei oleks osatud tähelepanu pöörata.

26 Teiselt poolt lihtsustab see tööd, kui sama isik satub järgmisel korral muu arsti juurde.

27 Tähtis tuleb ka meeles pidada, et eestlaski patsientide tähta vester küsitlusele varasemate haiguste ja operatsioonide, allergiate, põetud lastehaiguste, tarvitatavate ravimite, perekonnas esinenud haiguste ja muu kohta.

28 Eriti oluliseks teadmiseks on Eesti selliste küsimustike kasutamine eriti lühend pole.

29 Tähtselt Tartu Üldkoooli poliitilistest saadud informatsioonid põhjal soovitud neil selliste rühmid, kus nii

Plain Text Tab Width: 4 Ln 27, Col 79 INS

SISUKORD

Sissejuhatus 3

1. Sõnastiku sisse ja allikate valik 4

2. Eesti oskussõnastik 5

3. Patsiendi uurimine 6

3.1 Anamnees, läbivaatus ja instrumentaalne uuring 6

3.1.1 Küsitlus 7

3.1.2 Organismi talitlus ja elundite toimimine 8

3.1.3 Varasem tervislik seisund ja pöördumised 8

3.1.4 Sotsiaalsed olud ja eluviisid 8

3.1.5 Kliiniline uuring 8

3.2 Staatuse 9

3.3 Uurimisjärjekord 9

4. Haigusloog 10

4. Terminoloogia ja termin 10

5. Meditsiini oskussõnad 13

5.1 Eesti meditsiini terminoloogia 14

6. Sõnastiku koostamise tekkimise probleemid 15

7. Sõnastiku kasutamine 16

Haige uurimise saksakeelsete lühendid 18

Saksakeelset lühendit 32

Eesti-saksa registreeritud 33

Kokkuvõte 47

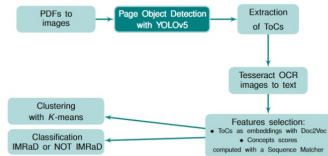
Kasutatud kirjandus 48

Resümee 55

Lisa 56

YOLO [6]

- Object detection algorithm
- Images divided into a grid. Objects detected within each grid cell
- Draws bounding boxes around regions of interest **AND** attribute them a label
- Used in our method: YOLOv5



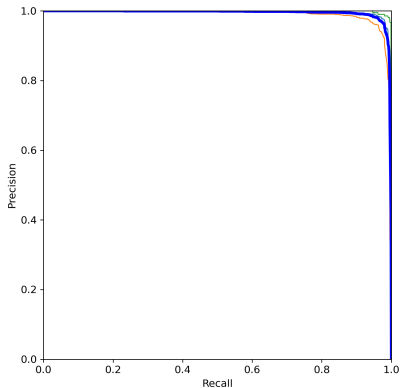
Training on our data

- 3 classes: 'headers', 'ToCs', 'body'
- Training set: 1400 images, testing set: 757 images
- Model: YOLOv5 small [6]
- Trained **from scratch** on GPU - 150 epochs

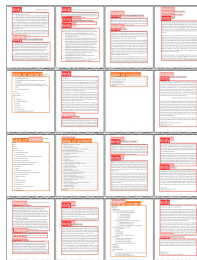
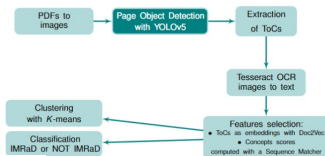
SISUKORD

SISUKORD	4
1. TEOHETILINE RAAMISTIK	7
1.1. Aia teoloogia ja teoloogiline teoloogia	7
1.2. Aia teoloogia teoloogia	12
1.3. Aia teoloogia teoloogia	17
1.4. Aia teoloogia teoloogia	20
1.5. Aia teoloogia teoloogia	23
2. MATERIAAL JA METOODIKA	25
2.1. INTERVJUUDE ANALÜÜS	32
2.1.1. Aia teoloogia teoloogia	32
2.1.2. Aia teoloogia teoloogia	38
2.1.3. Aia teoloogia teoloogia	40
2.1.4. Aia teoloogia teoloogia	41
2.1.5. Aia teoloogia teoloogia	42
4. VEEBIRKÜSTLUSTE TULEMISTE ANALÜÜS	45
4.1. Suhtumise erinevused	45
4.2. Vastajate arvamus aia teoloogia teoloogia	46
4.3. Aia teoloogia teoloogia	48
4.4. Aia teoloogia teoloogia	49
4.5. Aia teoloogia teoloogia	50
4.6. Aia teoloogia teoloogia	53
5. ARUTELU	58
5.1. Aia teoloogia teoloogia	58

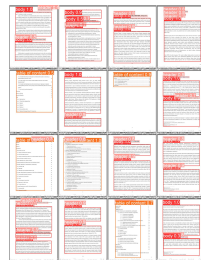
Results YOLOv5



Precision-Recall curve



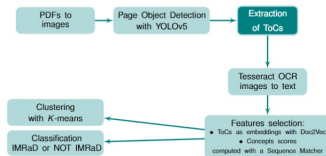
Labels



Predictions

Extraction: bounding boxes

- Extract of bounding boxes with coordinates from YOLOv5



Uusimäntägründit: koostajalike toimetetel ootikolhoon luteatun 2

Sisukord

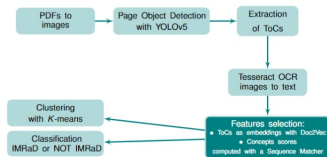
Sisukord	3
1. Toimetetelid lühikohad	4
1.1. Uusimäntägründit lühikohad	4
1.2. Uusimäntägründit lühikohad	7
1.3. Uusimäntägründit lühikohad	8
1.4. Uusimäntägründit lühikohad	12
2. Uusimäntägründit ja uusimäntägründit	15
3. Mõisted	16
3.1. Uusimäntägründit	16
3.2. Uusimäntägründit	17
3.3. Uusimäntägründit	17
4. Tulemused ja arutlus	19
4.1. Tulemused ja arutlus	19
4.2. Tulemused ja arutlus	22
4.3. Tulemused ja arutlus	28
4.4. Tulemused ja arutlus	41
4.5. Tulemused ja arutlus	43
Kokkuvõte	44
Summary	46
Tänuvõtte	48
Autoriteet	48
Kasutatud kirjandus	49
Lis 1	
Lis 2	
Lis 3	
Lis 4	
Lis 5	



Sisukord	3
1. Toimetetelid lühikohad	4
1.1. Uusimäntägründit lühikohad	4
1.2. Uusimäntägründit lühikohad	7
1.3. Uusimäntägründit lühikohad	8
1.4. Uusimäntägründit lühikohad	12
2. Uusimäntägründit ja uusimäntägründit	15
3. Mõisted	16
3.1. Uusimäntägründit	16
3.2. Uusimäntägründit	17
3.3. Uusimäntägründit	17
4. Tulemused ja arutlus	19
4.1. Tulemused ja arutlus	19
4.2. Tulemused ja arutlus	22
4.3. Tulemused ja arutlus	28
4.4. Tulemused ja arutlus	41
4.5. Tulemused ja arutlus	43
Kokkuvõte	44
Summary	46
Tänuvõtte	48
Autoriteet	48
Kasutatud kirjandus	49
Lis 1	

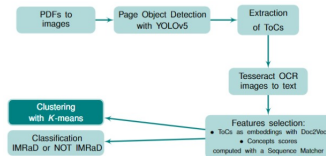
Feature selection

- ToCs → vectors of numerical values with Doc2Vec.
- Sequence Matcher [8]: compares and identifies similarities in given pairs of input strings.
- IMRaD concepts → concept scores.



<i>IMRaD</i> concepts	<i>IMRaD</i> words in Baltic languages
Introduction	'sissejuhatus', 'ivada', 'izanga'
Literature review	'valdkonna ulevaade', 'kirjanduse ulevaade', 'teoreetiline ulevaade', 'teoreetiline taust', 'teoreetiline raamistik', 'literatuuros apzvalga'
Methods	'metoodika', 'meetod', 'metodai', 'metodologija', 'metodika', 'metodine', 'metodas'
Results	'tulemused', 'uurimistulemused', 'rezultatai', 'duomenu analize ir rezultatu apzvalga'
Discussion	'arutelu', 'diskussioon', 'jareldused', 'aptarimas', 'isvados'

Clustering



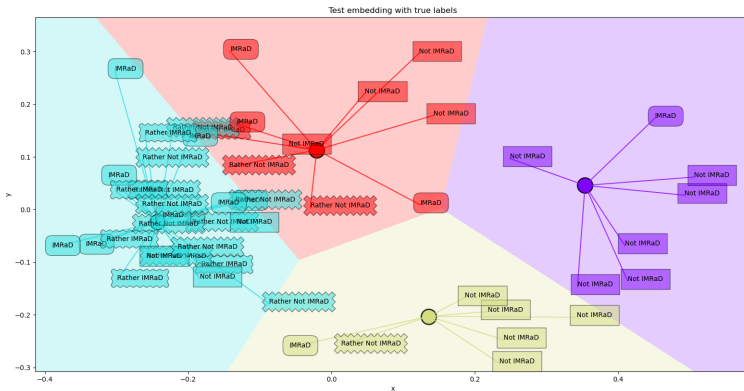
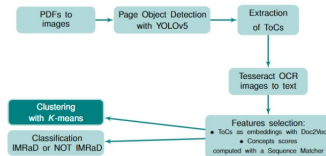
K-means clustering [4]

- Unsupervised method which groups similar data instances together
- Measures features selection's relevance
- 467 documents organized in four classes:
 - 98 *IMRaD* files
 - 171 *Not IMRaD* files
 - 105 *Rather IMRaD* files
 - 93 *Rather Not IMRaD* files
- Tested with $k = 4$ and $k = 2$
- Score embeddings, Text embeddings, Text and score embeddings

Clustering results

2D cluster representation

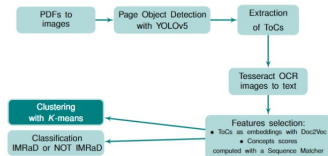
$k = 4$



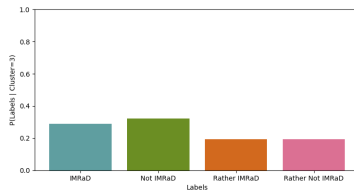
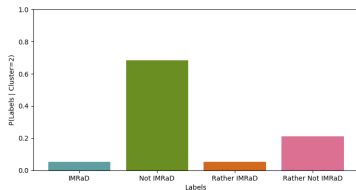
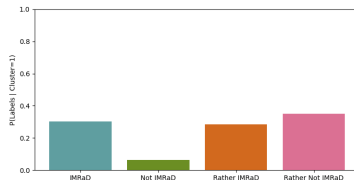
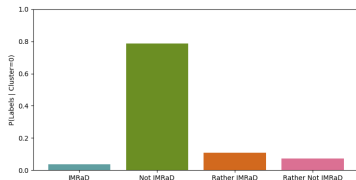
Clustering results

Histogram representation of distribution of classes across clusters.

$k = 4$



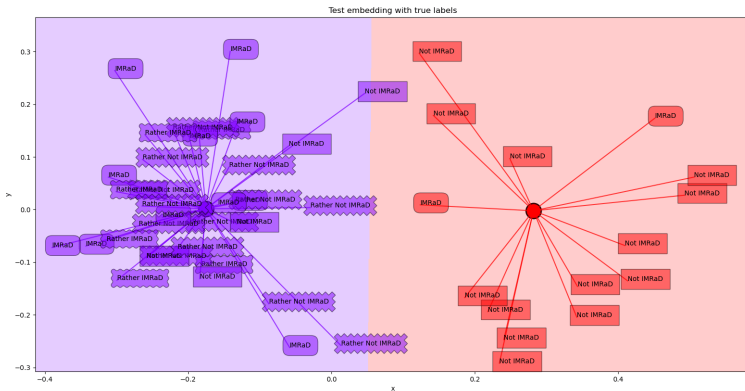
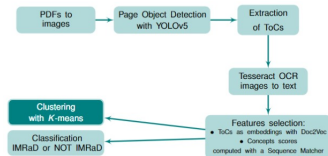
P(Labels | Clusters) on Test data



Clustering results

2D cluster representation

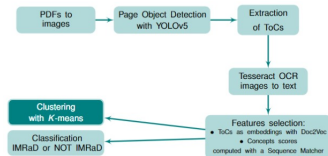
$k = 2$



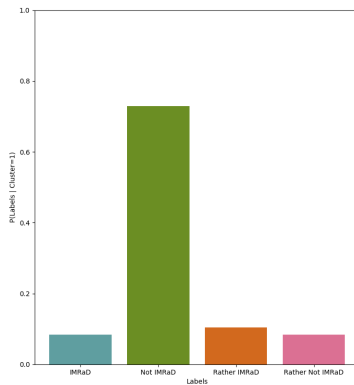
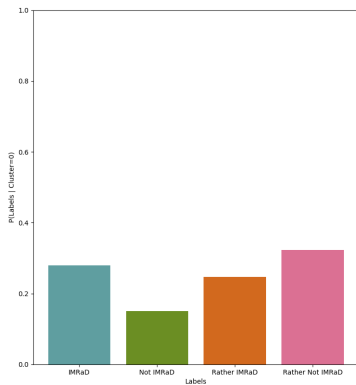
Clustering results

Histogram representation of distribution of classes across clusters.

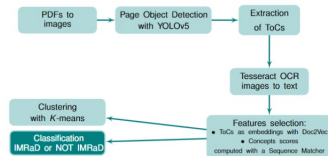
$k = 2$



P(Labels | Clusters) on Test data



Classification

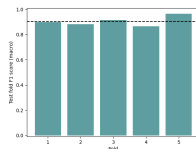
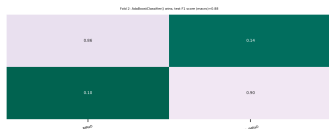
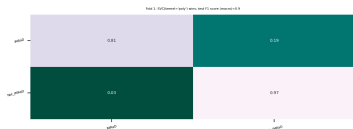
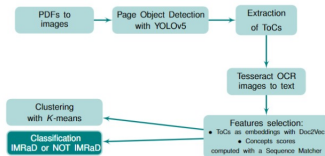


Def IMRaD			Def Not IMRaD
Def IMRaD	Rather IMRaD	Rather Not IMRaD	Def Not IMRaD
Def IMRaD	Rather IMRaD + Rather Not IMRaD		Def Not IMRaD
Def IMRaD + Rather IMRaD		Def Not IMRaD + Rather Not IMRaD	
Def IMRaD + Rather IMRaD + Rather Not IMRaD			Def Not IMRaD

Model selection

- 6 different models
- Variation in hyperparameters → total model tested: 32

Results classification



- Mean macro f1 score: 90%
- Standard deviation: 0.02

Discussion

- Pros and cons of object detection methods:
 - Usable with PDF documents without losing important graphical information and helps reduce the amount of data needing to be processed
 - Identification of headers in text remains challenging
- The possibility of detecting the organizational structure of a document independently of its language
- IMRaD
 - Not binary to a human annotator,
 - *Rather IMRaD* and *Rather Not IMRaD* documents not different enough from *IMRaD* files to a machine learning algorithm

References I



Kaplan, Robert B.

‘Cultural thought patterns in inter-cultural education’.
Language learning, 16(1-2):1–20, 1966.



Redmon, J., Divvala, R., Girshick, R., Farhadi, A.

‘You only look once: Unified, real-time object detection’.
Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 779–788, 2016.



Bochkovskiy, A. and Wang, C.-Y. and Liao, H.-Y. M.

‘Yolov4: Optimal speed and accuracy of object detection’.
arXiv preprint arXiv:2004.10934, 2020.



Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B.,
Grisel, O., and Duchesnay, E.

‘Scikit-learn: Machine learning in Python.’
the Journal of machine Learning research, 12, 2825–2830, 2011.

References II



Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel, O., and Varoquaux, G.

‘API design for machine learning software: experiences from the scikit-learn project.’

arXiv preprint arXiv:1309.0238., 2013.



Jocher, G., Nishimura, K., Mineeva, T. and Vilariño, R.

YOLOv5, 2020

<https://github.com/ultralytics/yolov5>



Hoffstaetter, S., Bochi, J., Lee, M., Kistner, L., Mitchell, R., Cecchini, E., Hagen, J., Morawiec, D., Bedada, E., and Akyüz, U.

Pytesseract

<https://github.com/madmaze/pytesseract>

References III



Python Software Foundation

Difflib - SequenceMatcher

<https://github.com/python/cpython/blob/master/Lib/difflib.py>



Cartucho, J., Ventura, R., Veloso, M.

Robust Object Recognition Through Symbiotic Deep Learning In Mobile Robots

In 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2336–2341, 2018.