

# ARGUMENT MINING ON CLINICAL TRIALS

TOBIAS MAYER

SUPERVISOR: SERENA VILLATA

CO-DIRECTOR: CÉLINE POUDAT

UNIVERSITÉ CÔTE D'AZUR, CNRS, INRIA, I3S  
BASE, CORPUS, LANGAGE (BCL)



UNIVERSITÉ  
CÔTE D'AZUR

## **Evidence-based medicine (EBM):**

- optimize decision making with evidence from well-conducted research
- meta-analysis and systematic reviews on Randomized Controlled Trials (RCT)

## **Evidence-based medicine (EBM):**

- optimize decision making with evidence from well-conducted research
- meta-analysis and systematic reviews on Randomized Controlled Trials (RCT)

How to assist with automatic processing?

## **Argument mining system for clinical trials:**

- automated approach to extract argumentative information from trials
- detection of claims and evidence
- domain unspecific applicability

"The general task of analyzing discourse on the pragmatics level and applying a certain argumentation theory to model and automatically analyze the data at hand"

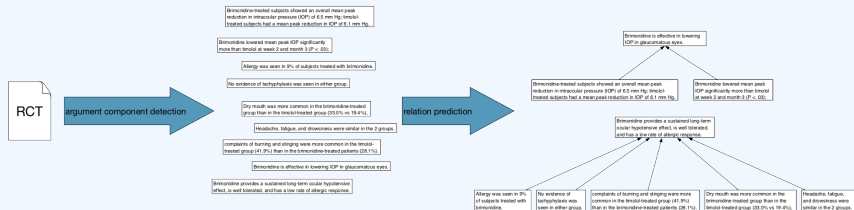
[Habernal and Gurevych, 2017]

## Argument extraction:

- distinguish argumentative from non-argumentative components
- classify the components into evidence and claims

## Relations prediction:

- intra-argument relation prediction
- inter-argument relation prediction



# DATASET

## Randomized Controlled Trials (RCT):

- common type of **experimental studies** in the medical domain
- **comparison** between intervention and control arm
- used for **evidence-based medical decision making** (systematic reviews and meta-analysis)
- **PubMed**: freely available citation database from the United States National Library of Medicine (NLM)
- structure should follow CONSORT<sup>1</sup> policies

---

<sup>1</sup><http://www.consort-statement.org/>



## Data collection:

- Annotate existing **collection of RCT abstracts** on glaucoma treatments with argumentative labels
- Delete existing PICO<sup>2</sup> annotations
- Extending the collection with more RCT abstracts from PubMed (**glaucoma, diabetes, hepatitis and hypertension**)

---

<sup>2</sup>Annotation framework for: **P**opulation, **I**ntervention, **C**ontrol and **O**utcome

## Claim

- concluding statement made by the author about the outcome of the study:
  - ▶ *"Brimonidine is well tolerated and has a low rate of allergic response."*
- general description of the relation between intervention and control arm:
  - ▶ *"Trabeculectomy was more effective than viscocanalostomy in lowering IOP in glaucomatous eyes of white patients."*
- should logically follow from the described results

## Evidence/Premise

- observation in the study (side-effect or other measured outcome):
  - ▶ *"Allergy was seen in 9% of subjects treated with brimonidine."*
  - ▶ *"Brimonidine lowered mean peak IOP significantly more than timolol at week 2 ( $P < .03$ )."*
- credible without further evidence (ground truth)
- supports or attacks another argument component

## ANNOTATION EXAMPLE

To compare the intraocular pressure-lowering effect of latanoprost with that of dorzolamide when added to timolol. [...] [The diurnal intraocular pressure reduction was significant in both groups ( $P < 0.001$ )]<sub>1</sub>. [The mean intraocular pressure reduction from baseline was 32% for the latanoprost plus timolol group and 20% for the dorzolamide plus timolol group]<sub>2</sub>. [The least square estimate of the mean diurnal intraocular pressure reduction after 3 months was -7.06 mm Hg in the latanoprost plus timolol group and -4.44 mm Hg in the dorzolamide plus timolol group ( $P < 0.001$ )]<sub>3</sub>. Drugs administered in both treatment groups were well tolerated. This study clearly showed that **[the additive diurnal intraocular pressure-lowering effect of latanoprost is superior to that of dorzolamide in patients treated with timolol]**<sub>1</sub>.<sup>3</sup>

---

<sup>3</sup>**claims** are written in bold, evidence are underlined

<i>Topic</i>	<i>#abstracts</i>	<i>#evidence</i>	<i>#claims</i>	<i>#non arguments</i>
glaucoma	119	448	153	743
diabetes	20	84	41	112
hepatitis	20	105	22	121
hypertension	20	60	33	126

## Inter-annotator agreement<sup>4</sup>:

- argumentative components: 0.72
- claim/evidence distinction: 0.68

---

<sup>4</sup>agreement is given in Fleiss' kappa for three annotators

# EXPERIMENTAL SETTINGS

## MARGOT<sup>5</sup>:

- argument mining approach to **overcome genre-dependency**
- addresses **argument component detection**
- **cross-domain features** (word occurrences, sentence structure)
- trained on Wikipedia articles

---

<sup>5</sup>MARGOT: Mining Arguments from Text. <http://margot.disi.unibo.it>

## Model:

- SVM classifier for detection of claim/evidence
- SVM+HMM for detection of component boundaries

## Features:

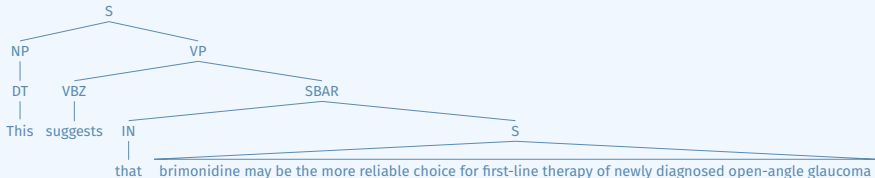
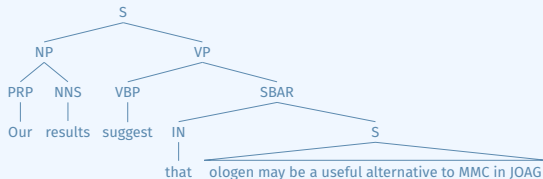
- SubSet Tree Kernels (SSTK)
- bag-of-words with TF-IDF values



## What is a tree kernel?

- **similarity measure** between constituency parse trees
- considers **common fragments** between two trees
- defines a rich feature space
- SSTK provides a **good compromise** between expressiveness and efficiency

# TREE KERNEL EXAMPLE



## Data partitioning:

topic	training	testing
glaucoma	79 abstracts	30 abstracts
hepatitis(HB)		20 abstracts
diabetes(DM)		20 abstracts
hypertension(HT)		20 abstracts
mixed		90 abstracts

# RESULTS

## RESULTS: EVIDENCE DETECTION<sup>6</sup>

		Glaucoma	DM	HB	HT	Mixed
<b>Evidence</b>	BoW	0.84	0.79	0.74	0.80	0.80
	SSTK	0.86	0.79	0.75	0.80	0.80
	SSTK + BoW	0.86	0.79	0.75	0.80	0.80

- SSTK slightly better than BoW, but still comparable
- no differences between SSTK and BoW for out-of-domain topics
- distinctive vocabulary might be related to statistical evaluation rather than medical terminology

---

<sup>6</sup>results are given in  $f_1$ -score

## RESULTS: CLAIM DETECTION<sup>6</sup>

		Glaucoma	DM	HB	HT	Mixed
Claim	BoW	0.75	0.68	0.62	0.64	0.65
	SSTK	0.79	0.73	0.66	0.70	0.72
	SSTK + BoW	0.79	0.74	0.66	0.70	0.72

- SSTK significantly better than BoW
- distinctive syntactic structure for claims
- SSTK generalizes better than BoW
- combining the models do not increase results
- lexical information also captured in syntactic structure

---

<sup>6</sup>results are given in  $f_1$ -score

## RESULTS: ARGUMENTATIVE COMPONENT DETECTION<sup>6</sup>

		Glaucoma	DM	HB	HT	Mixed
<b>Arg. Comp.</b>	BoW	0.82	0.74	0.70	0.72	0.74
	SSTK	0.86	0.76	0.71	0.74	0.78
	SSTK + BoW	0.86	0.76	0.71	0.74	0.78

- TK model performs better
- results similar to evidence detection
- many errors were made between claim and evidence distinction

---

<sup>6</sup>results are given in  $f_1$ -score

# EVIDENCE CLASSIFICATION



- EBM focuses mainly on study design and risk of bias as quality of evidence
- need for other aspects to measure trial quality (reproducibility, generalizability or estimate of effect)
- first step towards creating arguments for argumentation framework

## ■ **comparative:**

- ▶ *"The overall success rates were 87% for the 350-mm<sup>2</sup> group and 70% for the 500-mm<sup>2</sup> group ( $P = 0.05$ )."*

## ■ **significance:**

- ▶ *"All regimens produced clinically relevant and statistically significant ( $P < .05$ ) intraocular pressure reductions from baseline."*

## ■ **side-effect:**

- ▶ *"Allergy was seen in 9 % of subjects treated with brimonidine."*

## ■ **other:** risk factors, limitations

- ▶ *"Risk of all three outcomes was higher for participants with chronic kidney disease or frailty."*

Results for multi-class classification using SVMs:

<i>Dataset</i>	<i>Method</i>	<i>glaucoma</i>	<i>combined.</i>
Gold standard	RANDOM	0.33	0.32
	MAJORITY	0.27	0.26
	N-GRAMS	0.80	0.74
whole pipeline	RANDOM	0.38	0.38
	MAJORITY	0.38	0.39
	N-GRAMS	0.71	0.66

**Table:** Results (weighted average  $F_1$ -score).

- creation of a dataset of RCTs labeled with argumentative components
- application of Argument Mining on clinical trials
- first step to evidence classification

- relation prediction (building argumentation trees)
- annotation of CHU data (French) and corpus building (together with BCL)
- evidence quality assessment
- reproducible support for clinical decision making

THANK YOU FOR YOUR ATTENTION!

Description of the objective of a study confused as claims:

- *"The goal of this research is to evaluate efficacy and safety of herbal medicine as compared to allopathic medicine in patients suffering from hepatitis B."*

Claims sometimes with a very complex syntactic structure:

- *"The authors tested the hypothesis that a valsartan/cilnidipine combination would suppress the home morning blood pressure ( BP ) surge ( HMBPS ) more effectively than a valsartan/hydrochlorothiazide combination in patients with morning hypertension , defined as systolic BP ( SBP ) 135 mm Hg or diastolic BP 85 mm Hg assessed by a self-measuring information and communication technology-based home BP monitoring device more than three times before either combination 's administration."*



Group descriptions (group sizes or initial medical conditions) mis-classified as evidences:

- *"Among 426 participants (53% male, median age 35 years, median CD4 count 19 cells/ $\mu$ L), 31 developed hepatotoxicity (7.3%)."*
- *"Overall, there were no significant differences in pregnancy-induced hypertension across supplement groups."*

Negated sentences often mis-classified:

- *"No patients developed additional resistance mutations throughout the study period."*