# Improving Argument Mining in Student Essays by Learning and Exploiting Argument Indicators versus Essay Topics

Huy Nguyen[1]        Diane Litman[1,2]

[1]Computer Science Department
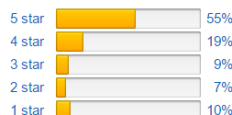
[2]Learning Research & Development Center

University of Pittsburgh, Pennsylvania

# Argumentative text is ubiquitous

"Argumentation mining […] involves automatically identifying **argumentative structures** within a document, […] as well as argument-subargument and argument-counterargument **relationships between pairs of arguments** in the document."

*(The 1st Argument Mining Workshop at NAACL, 2014)*

# Overview of our research

| |
|---|
| Essay evaluation |
| Argumentative relation classification |
| Argument component identification |
| Students' persuasive essays |

Do arts and music improve the quality of life?

My view is that the *government should give priorities to invest more money on the basic social welfares such as education and housing instead of subsidizing arts relative programs*MajorClaim.

Art is not the key determination of quality of life, but education isClaim. In order to make people better off, it is more urgent for governments to commit money to some fundamental help such as setting more scholarships in education section for all citizensPremise. This is simply because knowledge and wisdom is the guarantee of the enhancement of the quality of people's lives for a well-rounded social systemPremise.

Admittedly, art, to some extent, serve a valuable function about enriching one's daily livesClaim, for example, it could bring release one's heavy burden of study pressure and refresh human bodies through a hard day from workPremise. However, it is unrealistic to pursuit of this high standard of life in many developing countries, in which the basic housing supply has still been a huge problem with plenty of lower income family have squeezed in a small tight roomPremise. By comparison to these issues, the pursuit of art seems unimportant at allPremise.

To conclude, art could play an active role in improving the quality of people's livesPremise, but I think that governments should attach heavier weight to other social issues such as education and housing needsClaim because those are the most essential ways enable to make people a decent lifePremise.

# Argument component identification

- Argument component: text portion with a specific role in forming the argument[*]

[...] To conclude, art could play an active role in improving the quality of people's lives, but I think that governments should attach heavier weight to other social issues such as education and housing needs because those are the most essential ways enable to make people a decent life.

(Persuasive Essay Corpus, Stab & Gurevych 2014)

Premise 1

Attacks

Claim

Supports

Premise 2

- The step before argumentative relation mining

- This study focuses on argument component identification in student essays

# Prior argument component identification studies

- N-gram and production rule features (VP$\rightarrow$VBG NP) [Stab & Gurevych 2014]
  - ✗ Large and sparse feature space
  - ✗ Have not considered abstraction of argument topic

- Lexicons of argument and domain words [Nguyen & Litman 2015]
  - ✗ Lacked a quantitative evaluation

- Cross-fold validation
  - ✗ Have not evaluated topic-independence of the models (e.g., train and test essays are of different topics)

# Argument and domain word extraction [Nguyen & Litman 2015]

- 6794 un-annotated persuasive essays*

- Process Latent Dirichlet Allocation [Blei et al. 2003] **topic model output**

**Argument seeds**: agree, disagree, reason, support, advantage, disadvantage, think, conclusion, result, opinion

**Domain seeds**: *title words that are not argument seeds or stop words*

**Scoring algorithm**: *looks for the most argumentative LDA topic, i.e., high argument weight and low domain weight*

**Result**: *263 argument words and 1806 domain words (stemmed)*

believe

view

opinion

however

discuss

children

parent

school

learn

teach

music

art

creative

talent

idea

# Example argument and domain words

**Argument seeds**: agree, disagree, reason, support, advantage, disadvantage, think, conclusion, result, opinion

**LDA topic 1**: reason example support agree think because disagree statement opinion believe therefor idea conclusion

= *Argument seeds & variants, discourse connectives, stop words*

**LDA topic 2**: city live big house place area small apart town build community factory urban

**LDA topic 3**: children parent school education teach kid adult grow childhood behavior taught

# Baseline vs. Proposed models

**Stab14** (Stab & Gurevych 2014b)          **Nguyen15** (Nguyen & Litman 2015)          This study (**wLDA+4**)

*Lexical (I)*
- 1-, 2-, 3-grams
- Verbs, adverbs, presence of model verb
- Singular first person pronouns
- Discourse connectives

*(I)*
- Argument words as unigrams
- Same as Stab14

*Parse (II)*
- Production rules
- Tense of main verb
- #sub-clauses, depth of parse tree

*(II)*
- Argumentative subject-verb pairs
- Same as Stab14

*Structure (III)*
- #tokens, token ratio, #punctuation, sentence position, first/last paragraph, first/last sentence of paragraph

*(III)*

*Context (IV)*
- #tokens, #punctuation, #sub-clauses, modal verb in preceding/following sentences

*(IV)*
- Same as Stab14

Nguyen15 **+**
1. Numbers of common words with title and preceding sentence
2. Comparative & superlative adverbs and POS
3. Plural first person pronouns
4. Discourse relation labels

10-fold cross validation (data was randomly split into training and test sets)          Cross writing-prompt validation (training and test essays are of different prompts)

# Ablated models

- Replace argument and domain lexicons in wLDA+4 model

  - SEED model: uses only argument and domain seeds    *Extracted lexicons vs. Seed words*

  - woLDA model: does not use seed words or the two lexicons
    - Removes argument word features    *With lexicons vs. Without lexicon*
    - Uses all subject-verb pairs

# Persuasive Essay Corpus [Stab & Gurevych 2014]

- 90 persuasive essays*
  - MajorClaim
  - Claim
  - Premise

- 3 expert annotators
  - Accuracy 0.88
  - Krippendorff's $\alpha_U$ 0.72

> government should give priorities to invest more money on the basic social welfares [...]

> I think that governments should attach heavier weight to other social issues such as education and housing needs

> those are the most essential ways enable to make people a decent life

| MajorClaim | Claim | Premise | Non-argumentative |
|------------|-------|---------|-------------------|
| 90 | 429 | 1033 | 327 |

# Evaluation method

- Cross-fold validation
  - Randomly: 10-fold cross validation
  - By-prompt: cross writing prompt validation

- In each folding
  - Select top 100 features in training folds (InfoGain + Ranking)
  - Train prediction model with top 100 features
  - Record prediction output on the test fold

# Experimental results: cross validation

| | 10-fold cross validation | | | | | Cross-prompt validation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Stab14 | Nguyen15 | woLDA | SEED | wLDA+4 | Stab14 | Nguyen15 | woLDA | SEED | wLDA+4 |
| Accuracy | 0.787* | 0.792* | 0.780* | 0.781* | **0.805** | 0.780* | 0.796 | 0.774* | 0.776* | **0.807** |
| Kappa | 0.639* | 0.649* | 0.629* | 0.632* | **0.673** | 0.623* | 0.654+ | 0.618* | 0.623* | **0.675** |
| Precision | 0.741* | 0.745* | 0.746* | 0.740* | **0.763** | 0.722* | 0.757* | 0.751 | 0.734 | **0.771** |
| Recall | 0.694* | 0.698* | 0.695* | 0.695* | **0.720** | 0.670* | 0.695* | 0.681* | 0.686* | **0.722** |

Best values in bold. +: $p < 0.1$, *: $p < 0.05$ by T-test when comparing with wLDA+4

*Obtains comparable performances between two experiment settings*

*Proposed model (wLDA+4) performs the best in 10-fold cross validation*

12 groups:
- 11 single-prompt groups (73 essays)
- 1 mixed group of minor prompts (17 essays)

Prompts: school, technologies, prepared food …

# Experimental results: holdout test sets

| | Stab's test set | | Nguyen's test set | |
|---|---|---|---|---|
| | Stab's reported | wLDA+4 | Nguyen's reported | wLDA+4 |
| Accuracy | 0.77 | **0.82** | 0.83 | **0.84** |
| Kappa | – | **0.68** | 0.69 | **0.71** |
| F1 | 0.73 | **0.75** | 0.76 | **0.78** |
| Precision | 0.77 | **0.79** | 0.79 | **0.81** |
| Recall | 0.68 | **0.73** | 0.74 | **0.76** |

# Feature evaluation

- Among all top features used to train the models
  - 49% are argument words
  - 8% are argumentative subject-verb pairs

  *LDA-enabled features in Nguyen15*
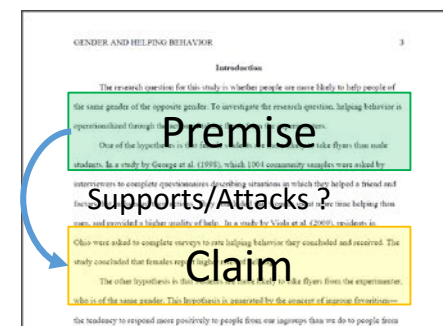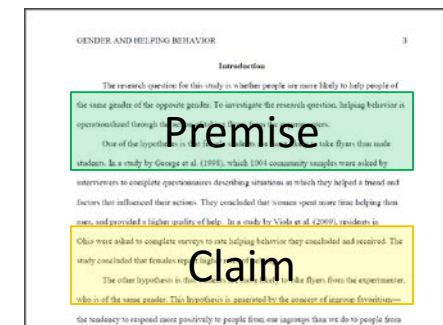
- In the top-50
  - Common word counts
  - Comparative adverbs, and *RBR* part-of-speech
  - Person pronouns *WE, OUR*
  - Discourse labels *Expansion, Contingency*

  *Proposed features in this study*

# Conclusions and future work



- A study on argument component identification

- New features that model argument indicators and abstract over essay topics
  - A necessary supplement to the learned and noisy argument and domain words

- Cross-topic and 10-fold cross validations
  - Proposed model obtained comparable performances



- Our next study focuses on argumentative relation classification

# Thank you!