

Homework 2 - Report

By

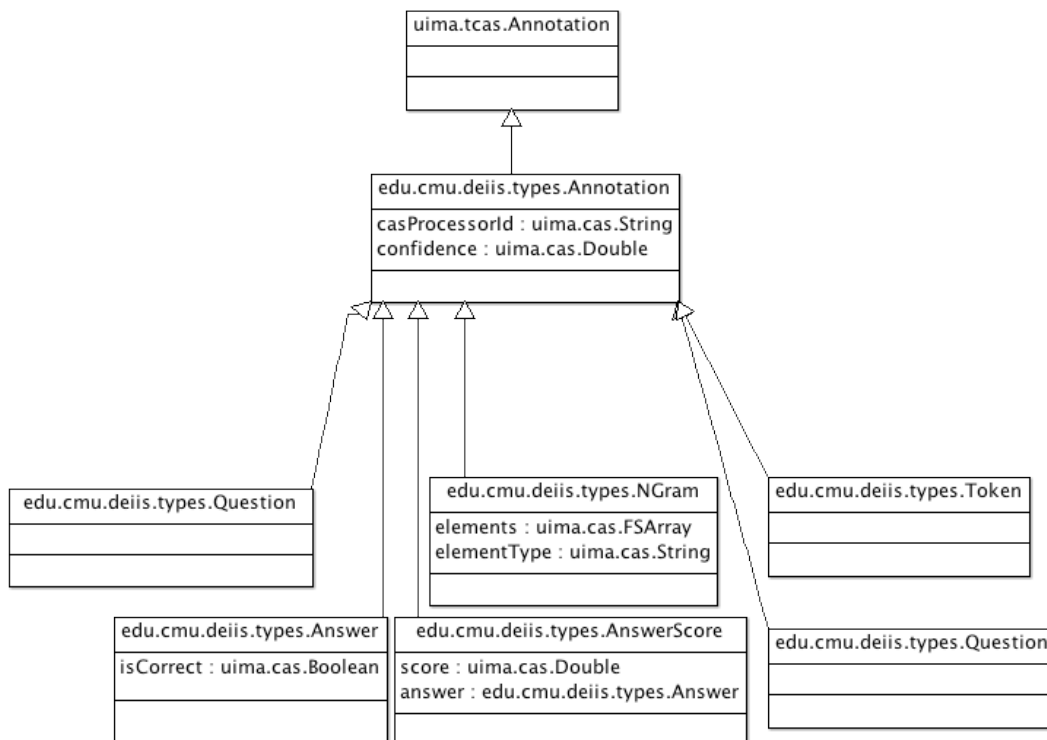
Vinay V Vemuri (Andrew ID: vvv)

Requirements

- 1. Annotate question from input text:** User provides an input file with question, proposed answers and a boolean flag indicating whether each answer is correct. System reads input file and annotates questions.
- 2. Annotate answers from input text:** User provides an input file with question, proposed answers and a boolean flag indicating whether each answer is correct. System reads input file and annotates all proposed answers. System also records whether each proposed answer is correct.
- 3. Annotate tokens in each question and answer:** The system takes in a sentence (question or an answer) as input and annotates all tokens in the sentence.
- 4. Annotate N-grams in each question and answer:** The system takes in a sentence (question or an answer) as input and annotates all 1-grams, 2-grams and 3-grams in the sentence.
- 5. Assign a score to each answer:** The system takes in an answer as input and assigns a score to the answer.
- 6. Sort answers according to scores and calculate precision:** The system will take in all answers and the score assigned to each answer as input and sort the answers according to their scores. Uses the total number of correct answers and the number of predicted correct answers to compute precision.

Given Type System

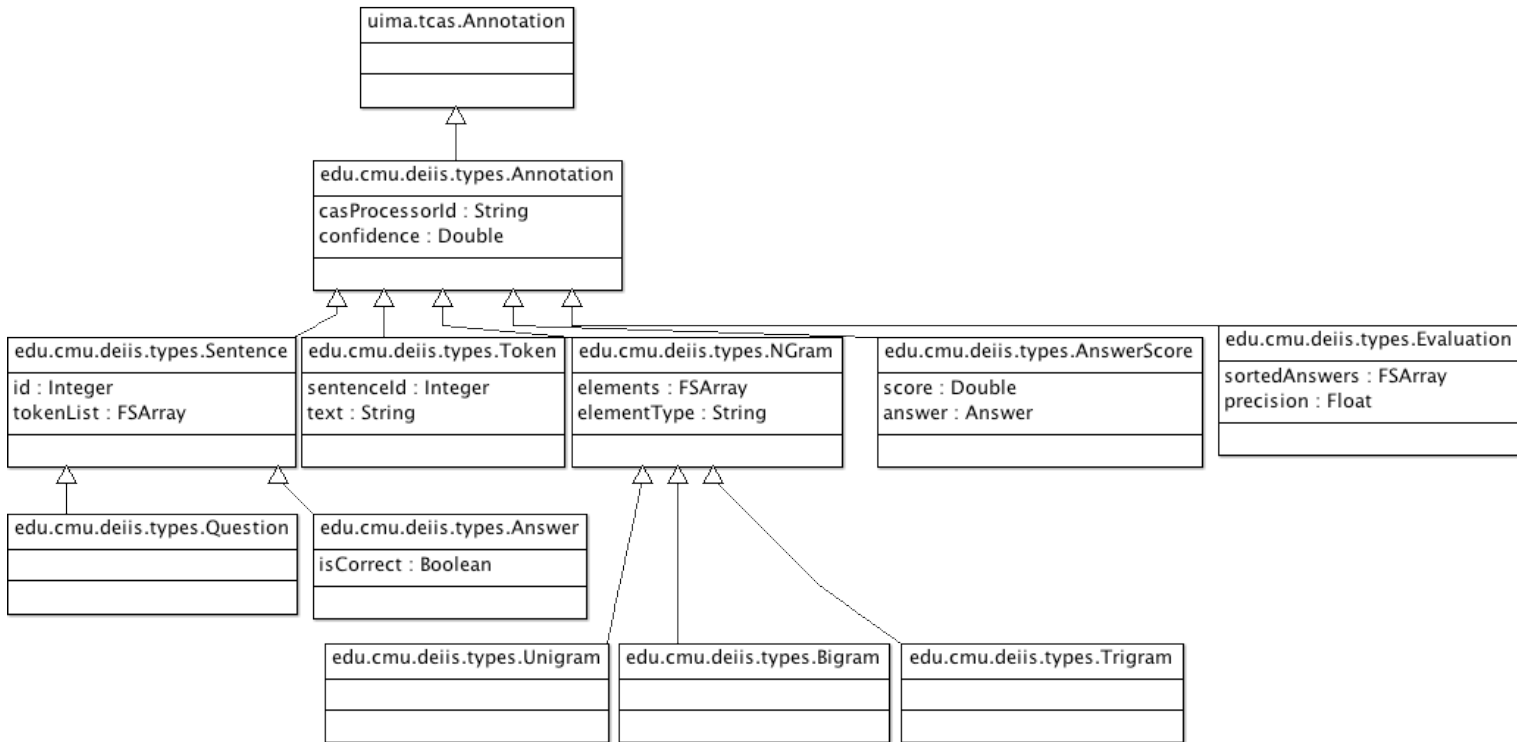
The UML Class diagram shown below provides a visualization of the given type system.



Homework 2 - Report

Improved Type-system

The following diagram shows the extended type system. Notice that no existing types from the given type system were deleted.



- **Annotation:** The Annotation type is the supertype of all the types in the proposed type system. The Annotation type inherits from the uima.tcas.Annotation type.

Its purpose is to ensure that all subtypes have a CasProcessorId (to indicate which class made the annotation) and a confidence value (to indicate how confident the source of the annotation was).

- **Sentence:** The Sentence type is the supertype of the Question and Answer types. It captures features common to both the Question and Answer types.

The 'Id' feature is a unique identifier assigned to each sentence in a given input file. The purpose of adding this feature was to allow each 'Token' to easily determine which sentence the Token is part of.

The TokenList feature of the sentence type is an array of all the Tokens in the sentence. The actual sentence text is captured as a feature of each Token in the TokenList.

- **Question & Answer:** The Question type inherits from the Sentence type. For the purposes of this project, the Question type does not have any more features other than the ones inherited from Sentence. However, more fields can be added as the need arises.

The Answer type also inherits from the Sentence type. However, unlike Question, the Answer has two other features in addition to the features inherited from Sentence. The 'IsCorrect' feature indicates whether the particular answer correctly answers the question.

Homework 2 - Report

- **AnswerScore:** The AnswerScore keeps track of the score assigned to the answer by the scoring system. It has two fields namely 'score' and 'answer'. The 'score' feature records the score assigned to the answer to allow answers to be sorted according to score during evaluation. The 'answer' feature keeps track of the answer for which the corresponding score has been recorded.
- **Token:** The Token type captures information about tokens in a Sentence (delimited by space and punctuation). It has two features namely SentenceId and Text. As mentioned previously, the 'SentenceId' is a unique identifier assigned to each sentence (and kept track of by each Token) to allow a Token to determine which sentence it is part of. The 'Text' feature captures the text contained in a token.
- **NGram:** The NGram is the supertype of Unigram, Bigram and Trigram and captures information about continuous sequences of tokens in a Sentence (Question or Answer). It consists of two features namely a list of Tokens named 'elements' of length N contained in the NGram and the type of element being stored in the elements array.
- **Unigram, Bigram & Trigram:** Unigram, Bigram and Trigram inherit from the NGram type. For the purposes of this project, these types don't have any additional fields to the ones they inherit from the NGram type. However, more fields can be added as the need arises.
- **Evaluation:** The Evaluation type captures information necessary to measure the performance of the entire system. This type consists of an array of Answers sorted by score to allow users to see which answers from input file best answer the question. This type also consists of a precision value to allow users to see the ratio of the actual number of correct answers to the predicted number of correct answers.

Screenshots

The following are screenshots of the UIMA Document Analyzer output with FullAnalysisEngine.xml.

Question

The screenshot displays the UIMA Document Analyzer output for a question. The main window shows the text "Q John loves Mary?" and a list of five possible answers with their corresponding scores. The "Question" type is selected in the legend. The right pane shows the annotation details for the "Question" type, including its begin and end positions, confidence, and a list of tokens.

Annotation Results for q002.txt.xmi in /Users/vvwmuri1/Desktop/Output

Q John loves Mary?
A 1 John loves Mary with all his heart.
A 1 Mary is dearly loved by John.
A 0 Mary doesn't love John.
A 0 John doesn't love Mary.
A 1 John loves Mary.

Legend

- ☐ Answer
- ☐ AnswerScore
- ☐ DocumentAnno...
- ☐ Evaluation
- ☐ NGram
- ☒ Question
- ☐ Token

Click In Text to See Annotation Detail

- Annotations
 - Question
 - Question ("John loves Mary?")
 - begin = 2
 - end = 18
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TestElementAnnotator
 - confidence = 1.0
 - id = 0
 - tokenList = FSArray
 - tokenList = Token ("John")
 - begin = 2
 - end = 6
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TestElementAnnotator
 - confidence = 0.0
 - sentenceId = 1
 - text = John
 - tokenList = Token ("loves")
 - tokenList = Token ("Mary")

Select All Deselect All Hide Unselected

Homework 2 - Report

Answer

Annotation Results for q002.txt.xmi in /Users/vvmemuri1/Desktop/Output

Q John loves Mary?
A 1 John loves Mary with all his heart.
A 1 Mary is dearly loved by John.
A 0 Mary doesn't love John.
A 0 John doesn't love Mary.
A 1 John loves Mary.

Click In Text to See Annotation Detail

- Annotations
 - Answer
 - Answer ("John loves Mary with all his heart.")
 - begin = 24
 - end = 59
 - casProcessorId = edu.cmu.deiis.AnnotationImpl.TestElementAnnotator
 - confidence = 1.0
 - id = 1
 - tokenList = FSArray
 - tokenList = Token ("John")
 - begin = 24
 - end = 28
 - casProcessorId = edu.cmu.deiis.AnnotationImpl.TestElementAnnotator
 - confidence = 0.0
 - sentencelid = 2
 - text = John
 - tokenList = Token ("loves")
 - tokenList = Token ("Mary")
 - tokenList = Token ("with")
 - tokenList = Token ("all")
 - tokenList = Token ("his")
 - tokenList = Token ("heart.")
 - isCorrect = true

Legend

☒ Answer ☐ AnswerScore ☐ DocumentAnn... ☐ Evaluation ☐ NGram

☐ Question ☐ Token

Select All Deselect All Hide Unselected

Token

Annotation Results for q002.txt.xmi in /Users/vvmemuri1/Desktop/Output

Q John loves Mary?
A 1 John loves Mary with all his heart.
A 1 Mary is dearly loved by John.
A 0 Mary doesn't love John.
A 0 John doesn't love Mary.
A 1 John loves Mary.

Click In Text to See Annotation Detail

- Annotations
 - Token
 - Token ("loves")
 - begin = 7
 - end = 12
 - casProcessorId = edu.cmu.deiis.AnnotationImpl.TokenAnnotator
 - confidence = 1.0
 - sentencelid = 0
 - text = loves

Legend

☐ Answer ☐ AnswerScore ☐ DocumentAnn... ☐ Evaluation ☐ NGram

☐ Question ☒ Token

Select All Deselect All Hide Unselected

Homework 2 - Report

NGram

Annotation Results for q002.txt.xmi in /Users/vvemuri1/Desktop/Output

Q John loves Mary?
A 1 John loves Mary with all his heart.
A 1 Mary is dearly loved by John.
A 0 Mary doesn't love John.
A 0 John doesn't love Mary.
A 1 John loves Mary.

Click In Text to See Annotation Detail

Annotations

- NGram
 - NGram ("dearly")
 - NGram ("dearly loved")
 - NGram ("dearly loved by")
 - NGram ("is dearly")
 - NGram ("is dearly loved")
 - NGram ("Mary is dearly")
 - begin = 65
 - end = 79
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.NGramAnnotator
 - confidence = 1.0
 - elements = FSArray
 - elements = Token ("Mary")
 - begin = 65
 - end = 69
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TokenAnnotator
 - confidence = 1.0
 - sentencId = 2
 - text = Mary
 - elements = Token ("is")
 - elements = Token ("dearly")
 - elementType = edu.cmu.deiis.types.Token

Legend

☐ Answer ☐ AnswerScore ☐ DocumentAnn... ☐ Evaluation ☒ NGram
☐ Question ☐ Token

Select All Deselect All Hide Unselected

AnswerScore

Annotation Results for q002.txt.xmi in /Users/vvemuri1/Desktop/Output

Q John loves Mary?
A 1 John loves Mary with all his heart.
A 1 Mary is dearly loved by John.
A 0 Mary doesn't love John.
A 0 John doesn't love Mary.
A 1 John loves Mary.

Click In Text to See Annotation Detail

Annotations

- AnswerScore
 - AnswerScore ("John loves Mary with all his heart.")
 - begin = 24
 - end = 59
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.AnswerScoringAnnotator
 - confidence = 0.0
 - score = 1.0
 - answer = Answer ("John loves Mary with all his heart.")
 - begin = 24
 - end = 59
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TestElementAnnotator
 - confidence = 1.0
 - id = 1
 - tokenList = FSArray
 - tokenList = Token ("John")
 - begin = 24
 - end = 28
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TestElementAnnotator
 - confidence = 0.0
 - sentencId = 2
 - text = John
 - tokenList = Token ("loves")
 - tokenList = Token ("Mary")
 - tokenList = Token ("with")
 - tokenList = Token ("all")
 - tokenList = Token ("his")
 - tokenList = Token ("heart.")
 - isCorrect = true

Legend

☐ Answer ☒ AnswerScore ☐ DocumentAnn... ☐ Evaluation ☐ NGram
☐ Question ☐ Token

Select All Deselect All Hide Unselected

Homework 2 - Report

Evaluation

Annotation Results for q002.txt.xmi in /Users/vvmemuri1/Desktop/Output

Q John loves Mary?
A 1 John loves Mary with all his heart.
A 1 Mary is dearly loved by John.
A 0 Mary doesn't love John.
A 0 John doesn't love Mary.
A 1 John loves Mary.

Legend

<input type="checkbox"/> Answer	<input type="checkbox"/> AnswerScore	<input type="checkbox"/> DocumentAn...	<input checked="" type="checkbox"/> Evaluation	<input type="checkbox"/> NGram
<input type="checkbox"/> Question	<input type="checkbox"/> Token			

Select All Deselect All Hide Unselected

Click In Text to See Annotation Detail

Annotations

- ▼ Evaluation
 - ▼ Evaluation ("John loves Mary with all his heart. A 1 Mary is dearly loved by...")
 - begin = 24
 - end = 174
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.EvaluationAnnotator
 - confidence = 1.0
 - ▼ sortedAnswers = FSArray
 - ▼ sortedAnswers = Answer ("John loves Mary with all his heart.")
 - begin = 24
 - end = 59
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TestElementAnnotator
 - confidence = 1.0
 - id = 1
 - ▼ tokenList = FSArray
 - ▼ tokenList = Token ("John")
 - begin = 24
 - end = 28
 - casProcessorId = edu.cmu.deiis.AnnotatorsImpl.TestElementAnnotator
 - confidence = 0.0
 - sentencId = 2
 - text = John
 - ▶ tokenList = Token ("loves")
 - ▶ tokenList = Token ("Mary")
 - ▶ tokenList = Token ("with")
 - ▶ tokenList = Token ("all")
 - ▶ tokenList = Token ("his")
 - ▶ tokenList = Token ("heart.")
 - isCorrect = true
 - ▶ sortedAnswers = Answer ("Mary is dearly loved by John.")
 - ▶ sortedAnswers = Answer ("John loves Mary.")
 - ▶ sortedAnswers = Answer ("Mary doesn't love John.")
 - ▶ sortedAnswers = Answer ("John doesn't love Mary.")
 - precision = 0.6