



# CMPE 493 Term Project Final Presentation

Named Entity Normalization for the Bacteria Biomes Domain

Şadi Uysal, Enes Özcan, Sercan Ersoy



# Previous Work

- Exact Match
  - Text exactly matches the name of an ontology term
- Jaccard Similarity
  - Jaccard similarity btw every ontology term's name
  - The term with the maximum score is predicted
- Cosine Similarity
  - Cosine similarity btw every ontology term's name
  - The term with the maximum score is predicted



# Previous Results

- Scores for the train data
- Total of 1118 habitats

	<b>Success</b>	<b>Fail</b>	<b>No Prediction</b>	<b>Accuracy (predicted)</b>	<b>Accuracy (total)</b>
<b>Exact Match</b>	166	43	909	79,42	14,85
<b>Jaccard Sim</b>	323	795	0	28,89	28,89
<b>Cosine Sim</b>	341	777	0	30,50	30,50



# New Prediction Logic

- Use exact matching result as a feature for score calculation.(1 or 0)
- Use exact matching **among synonyms** result as a feature for score calculation.(1 or 0)
- Use jaccard similarity result as a feature for score calculation.([0,1])



# New Prediction Logic

- IF calculated score is greater than threshold  
-return that candidate
- ELSE calculate **max** cosine similarity



# (1) Exact Matching

The same matching logic with the progress presentation is used.

- For each `'Habitat'` annotation in the *a1* files of the train data
- If annotated text exactly matches a term's name in the ontology (*obo* file)
- Then predicts that term as referent to the annotation
- Example:
  - In *a1*: `('T4', 'Habitat', 'gastric mucosa')`
  - In ontology: `'OBT:001792' -> 'gastric mucosa'`
  - Predicts `'OBT:001792'` as referent term



## (2) Exact Similarity Matching

Similar to exact matching, but among synonyms of an entity.

- For each 'Habitat' annotation in the *a1* files of the train data
- If annotated text exactly matches one of the synonyms of this term (including EXACT and RELATED ones)
- Then predicts that term as referent to the annotation
- Example:
  - In *a1*: ('T4', 'Habitat', 'Neisser +')
  - In ontology: id: OBT:000473
  - name: Neisser-positive
  - synonym: "Neisser +" RELATED [TyDI:57816]
  - is\_a: OBT:000174 ! Neisser stain phenotype
  - Predicts 'OBT:000473' as referent term



## (3) Weighted Similarity Matching

Combination of 3 features that we used.

- Exact match, exact similarity match, jaccard similarity
- `Exact_match=0/1` `Exact_match_synonyms=0/1` `Jaccard_coef=[0-1]`
- We assigned weights and calculated scores based on weights
- Hyper-parameter optimization





# Bidirectional Encoder Representations from Transformers

- BERT: Neural network-based technique for natural language processing pre-training
- huggingface/transformers library provides state-of-the-art general-purpose architectures for NLP with over thousands of pretrained models in 100+ languages
- BioBERT : A language representation model for biomedical domain



# Dictionaries

- ID\_Vector\_Dict\_Name={}


ID-vector dict for representations from id's\_name

- ID\_Vector\_Dict\_Training\_Data={}

ID-vector dict for representations from id's training data

- ID\_Vector\_Dict\_Ontology\_Classes={}

ID-vector dict for representations from id's ontology super classes

- 
- If we can not find candidate with  $\text{weighted similarity} > \text{threshold}$  :
    - Find id with max cosine similarity (with a minimum cosine similarity threshold)
- by using our dictionaries
- Return id



# Results

Development Set

<b>Habitats</b>	<b>0.4727</b>
<b>Habitats (exact)</b>	<b>0.3529</b>
<b>Habitats (new in dev)</b>	<b>0.4530</b>
<b>Habitats (new in test)</b>	<b>-</b>
<b>Habitats (only unique form-normalization)</b>	<b>0.3897</b>

Test Set

<b>Habitats</b>	<b>0.4500</b>
<b>Habitats (exact)</b>	<b>0.3147</b>
<b>Habitats (new in dev)</b>	<b>0.2727</b>
<b>Habitats (new in test)</b>	<b>0.1429</b>
<b>Habitats (only unique form-normalization)</b>	<b>0.4079</b>



# Things to improve

- Syntactic dependency parser, getting headwords and using those relations



# Thanks for your attention

Şadi Uysal

Enes Özcan

Sercan Ersoy