

# Problem Statement

<b>Title</b>	<b>Inconsistency Detection in Medical Annotation</b>
<b>Description</b>	<p>Text in clinical documents is annotated either automatically or manually to tag the words or phrases as problems, diagnoses, procedures, lab-tests, drugs, body parts etc and identify the relationship between them. Semantically similar words or phrases can be annotated differently by different annotators or maybe same annotator over time due to some ambiguity in understanding the guidelines, incomplete experience or human errors.</p> <p>Identifying these inconsistencies can help in improving the data which in turn helps in developing more effective models / learning by machine.</p> <p><b>Task</b> : Develop a tool which finds and display these similar entities which have different patterns of annotation in the supplied corpus.</p>
<b>YouTube Link</b>	<a href="https://www.youtube.com/watch?v=vx4PNSpE1wg">https://www.youtube.com/watch?v=vx4PNSpE1wg</a>
<b>Nature/Difficulty</b>	Complex
<b>Category:</b>	Software - Healthcare
<b>Technology Bucket:</b>	Software - Natural Language Processing, Software Development
<b>Dataset</b>	Will be provided in Hackathon
<b>For more details, contact</b>	Vivek Kumar ( <a href="mailto:vivek.k@ezdi.us">vivek.k@ezdi.us</a> )

# Detailed Description

## Introduction:

Annotation enables us to add information to plain text. Clinical domain is vast and diverse in nature and annotation of medical entities helps us to convert unstructured data to its structured format. Text of medical domain is annotated for clinical named entities with categories like problem, diagnosis, procedure, lab-test, medicine, anatomical structure etc. The relationship between mutually related entities is annotated as well. These annotations can be on single word or multi-word clinical entities.

### Example:

- The patient has left leg pain.
  - *'left'* -> 'modifier'
  - *'leg'* -> 'anatomical structure'
  - *'pain'* -> 'problem'

Similar words in a multi-word setting can be annotated differently by different annotators or maybe same annotator over time due to some ambiguity inherent in language or ambiguity in understanding the guidelines, incomplete experience or human errors. Two words or word sequences (Multi-word entities **MWEs**) which represent semantically similar category of meaning, if annotated differently, in terms of annotation label or boundary of entity, are considered inconsistently annotated.

**Example 1:** '*High blood pressure*' and 'Hypertension' are synonymous, but can be tagged in following ways:

- *'Hypertension'* -> 'Problem'
- *'high blood pressure'* -> 'Problem'
- *'high'* -> 'Modifier'      *'blood pressure'* -> 'Body measurement'

**Example 2:** '*Acute Kidney Injury*' can be annotated as:

- *'Acute Kidney Injury'* -> 'Problem'
- *'Acute'* -> 'Modifier'      *'Kidney Injury'* -> 'Problem'
- *'Acute'* -> 'Modifier'      *'Kidney'* -> 'Anatomical structure'      *'Injury'* -> 'Problem'

### Example 3:

1. **Internal fixation** will be done for the fracture.
  - *'internal fixation'* -> 'PROCEDURE'

2. **Fracture osteosynthesis** will be performed.
  - '*Fracture osteosynthesis*' -> 'PROCEDURE' (Synonym for "Internal Fixation")
  - '*Fracture*' -> 'PROBLEM', '*osteosynthesis*' -> 'PROCEDURE'
3. The patient underwent **open reduction** and **internal fixation** of left ankle.
  - '*open reduction and internal fixation*' -> 'PROCEDURE' (Also known as ORIF, single entity)
  - '*open reduction*' -> 'PROCEDURE', '*internal fixation*' -> 'PROCEDURE' (Annotated as two different procedures)
4. The physician will give recommendations with regards to further **radiographic evaluation** and **therapeutic intervention**.
  - '*radiographic evaluation*' -> 'PROCEDURE', '*therapeutic intervention*' -> 'PROCEDURE' (Annotated as two different procedures)

For **Sentence-2**, The definition of **Osteosynthesis** covers the fracture part, so annotators may be confused as to take "*Fracture Osteosynthesis*" as single entity, or as two different entities. Now look at **Sentence-1**, where "*internal fixation*" is marked as "PROCEDURE", but in **Sentence-3**, when "*Open Reduction*" is present with "*Internal Fixation*", this is usually treated as a single entity, so in this case also, annotators may decide to use the first annotation and second annotation is removed. Looking at similar type of content in **Sentence-4**, still two procedures are given, but they are always considered as independent entities and never combined together.

Now once all the inconsistencies are taken care of and the annotators mutually agree upon an annotation scheme. The task of correcting the annotation from the found clusters (or groups), can be performed semi-automatically.

**Impact:** Identifying these inconsistencies can help in improving the data quality. Better data helps in:

- Developing consistent and specific annotation schemes / rules which multiple human annotators can use to generate data which is less ambiguous and mutual disagreements which may arise due to human perspective.
- Human annotators to look back previously annotated examples for making decision while annotating new data.
- Developing more efficient machine learning (ML) models. Even a good ML algorithm cannot derive a good model if the training data fed to it is inconsistent. An ML system trained with consistently annotated data should yield better performance than one trained with inconsistent annotation.

**Task :**

- Find the **inconsistencies** in annotation.
  - Group the similar entities into various clusters.
  - Each cluster of entities is associated with various patterns of annotation available in the given corpus.
- Develop an **application** which displays the contents of these clusters and existing annotation patterns for conflicts.

**Possible Approaches for finding inconsistencies:**

- Dictionary lookup based - UMLS Lookup for individual entities, similar patterns and then find inconsistencies.
- Machine learning based - Grouping similar entities from sentences and then find conflicting patterns for those entities.

**Description of Annotations :**

- **Explanation of Annotation Tags** - The corpus is annotated into following 11 entity types:-

Entity Type	Description	Example
Problem (DG)	The disease conditions which include major problem, disease, symptoms & disorders.	<i>Complication of <b>bleeding, infection, arterial puncture, DVT.</b></i>
Finding	Concepts apart from major problem including abnormal conditions and minor alteration to the regular condition.	<i>This is a 27-year-old Caucasian female <b>gravida 4 para 1</b>, feeling <b>weak and lethargic.</b></i>
Procedure (PROC)	Surgery or other procedures performed for cure or diagnosis.	<i>This is an 82-year-old female with history of <b>appendectomy</b>, status post <b>open reduction internal fixation.</b></i>
Anatomical Structure	Anatomical sites, cells and organs of human body.	<i>The patient continued to have mild colitis throughout into the <b>cecum.</b></i>
Body Function	Activities carried out by the body to maintain the normal functioning of the body.	<i>The patient's <b>breathing</b> was normal.</i>
Lab Data	The type of analysis performed on blood, urine, other body substances or tissues to help diagnose or monitor the patient's condition.	<i><b>AFB testing</b> is negative, <b>TSH</b> is normal at 1.1.</i>

Body Measurement	The normal measurement of the body obtained without performing complex procedure or test.	<i>The <b>weight</b> is up a couple pounds at 157 pounds, <b>pulse</b> is 74.</i>
Measurement Value	Numerical value with its unit, associated with body measurement.	<i>The patient's heart rate was <b>90</b> and the blood pressure was <b>140/90</b>.</i>
Medical Device (MD)	Instruments used for treatment, operation and various medical purposes.	<i><b>Arterial line catheter</b> was placed over the <b>guidewire</b> without any resistance.</i>
Medicine	A drug used for the treatment or prevention of a disease.	<i>Completed antibiotic course of <b>ceftriaxone</b>.</i>
Modifier	Any word that adds some specific meaning to an Entity.	<i><b>Chronic</b> skin excoriation due to known neurodermatitis.</i>

**Table 1:** Entity types with description and examples.

## - Types of Relationships

### 1. Anatomical structure and problem relationship

This type of relationship helps to understand which part of the anatomical structure is affected by the particular problem. It simplifies the understanding of the problem and its area of effect.

**For example:** she had severe *pain* in the left *ankle*.

Here *pain* is the problem and *ankle* is the anatomical structure. Relationship can be formed between these two concepts and it makes it easy to understand that *pain* is in the *ankle*.

### 2. Anatomical structure and finding relationship

This relation helps to understand to which anatomical structure is the finding related.

**For example:** He noted he had felt some *tingling* and *numbness* in the left *upper chest* area the night before.

Here *tingling* is a finding and *chest* is the anatomical structure. By linking these two concepts in a relation we explain that tingling is occurring in the chest.

### 3. Anatomical structure and procedure relationship

This relation explains that a procedure is being done at given the anatomical structure. It is usually used to relate any operation, test, and other such procedures to the organ or body site at which it is being performed.

**For example:** The patient had a *CT* of the *brain*.

Here *CT* is a Computed Tomography procedure and *brain* is the anatomical structure. Relating these two concepts simplifies the understanding that *CT* of the *brain* was performed.

#### 4. Body measurement and Measurement Value relationship

Body measurement is the measurement of basic body parameters like temperature, height, weight, pulse rate, etc. These parameters often have values that can be linked to the measurement in order to provide a complete meaning.

**For example:** *Pulse: 80. BP: 110/70. Respirations: 16. Temp: 97.4*

Here, pulse, BP, temperature, and respiration are the body measurements and 80, 110/70, 16, and 97.4 are the measurement values of these respectively. So, forming the relationship between pulse and 80 signifies that pulse rate is 80 per min, and similarly for all such concepts.

#### 5. Relationships of Modifier

All the concepts except modifier itself and measurement values can be related to the modifier to add specific meaning to these concepts. Modifiers can never stand alone.

### Preparation of train data:

- **De-identification of Sentences** - As we are dealing with clinical data, de identification of this data is very important. Any personally identifiable information of a patient cannot be disclosed. The data will be deidentified to remove all personal information like names and case number, dates, age etc.
- **Dataset to be provided:**
  - **Size** - We will select a portion from annotated database, about 10,000 sentences.
  - **Type** - Each sentence will be of the format:
    - *He denies any issues this morning with [chest pain]**DG**, [worsening dyspnea]**DG**, [abdominal pain]**DG** , [nausea]**DG** , or [vomiting]**DG***
    - *He has supplemental oxygen in place via [nasal cannulae]**MD***
    - *[Drainage]**PROC** was carried out to relieve his [pain]**DG**.*
  - **DG** - Diagnosis/Problem, **MD** - Medical Device, **PROC** - Procedure

### Evaluation Criterion:

- We will prepare clusters of similar entities and associate the available annotation patterns of them.
- The participants' output also will be clusters of similar entities and their annotation patterns.
- A tool will compare the output clusters and gold clusters and calculate a similarity score.
- This similarity score will be the criterion of performance of the participants' tools.

## Supporting Material:

### Annotation Guidelines :

- [Annotation guidelines](#)
- Nancy Ide and Laurent Romary, "International Standards for Linguistic Annotation Framework" <https://arxiv.org/pdf/0707.3269.pdf>
- [https://www.researchgate.net/publication/318175802\\_Introduction\\_The\\_Handbook\\_of\\_Linguistic\\_Annotation](https://www.researchgate.net/publication/318175802_Introduction_The_Handbook_of_Linguistic_Annotation)
- 

**UMLS :** <https://www.nlm.nih.gov/research/umls/>

### Research Papers containing similar problems :

- Nora Hollenstein, Nathan Schneider, Bonnie Webber. 2016. Inconsistency Detection in Semantic Annotation, In Proc. LREC-2016 - [http://www.lrec-conf.org/proceedings/lrec2016/pdf/584\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2016/pdf/584_Paper.pdf)
- Markus Dickinson, Chong Min Lee. 2008. Detecting Errors in Semantic Annotation. In Proc. LREC-2008 - [http://www.cs.brandeis.edu/~marc/misc/proceedings/lrec-2008/pdf/157\\_paper.pdf](http://www.cs.brandeis.edu/~marc/misc/proceedings/lrec-2008/pdf/157_paper.pdf)