BIBM 2008 November 3-5, Philadelphia, PA

IEEE International Conference on
Bioinformatics and Biomedicine
Philadelphia, Pennsylvania
November 5, 2008

# Ontologies for Mining Biomedical Data

NATIONAL
LIBRARY OF
MEDICINE

*Olivier Bodenreider*

Lister Hill National Center
for Biomedical Communications
Bethesda, Maryland - USA

# Outline

- ◆ Biomedical ontologies
  - What they are
  - What they are for
  - Examples
- ◆ Ontologies for mining biomedical data
  - Normalization
  - Integration
  - Aggregation
- ◆ Applications of ontologies to data mining in biomedicine
  - Text mining – Information extraction
  - Biological – Mining gene expression data and functional annotations
  - Clinical – Mining adverse drug reactions

NLM

2

# Biomedical ontologies

# What is an ontology?

◆ The *What* question

- Objects in the world
  - With their properties
  - With their relations to other objects
- Also: events, processes, and states

◆ Explicit specification of a conceptualization

- Support software applications                    [Gruber, 1993]

◆ Domain ontology reflects

- Underlying reality
- Theory of the domain

4

# Ontology vs. other artifacts

- ◆ **Ontology**
  - Defining types of things and their relations
- ◆ **Terminology**
  - Naming things in a domain
- ◆ **Thesaurus**
  - Organizing things for a given purpose
- ◆ **Classification**
  - Placing things into (arbitrary) classes
- ◆ **Knowledge bases**
  - Assertional vs. definitional knowledge

# Examples of biomedical ontologies

◆ Structural perspective
- What are they (vs. what are they for)?

◆ "High-impact" biomedical ontologies
- International Classification of Diseases (ICD)
- Logical Observation Identifiers, Names and Codes (LOINC)
- SNOMED Clinical Terms
- Foundational Model of Anatomy
- Gene Ontology
- RxNorm
- Medical Subject Headings (MeSH)
- NCI Thesaurus
- Unified Medical Language System (UMLS)

NLM

6

# Characteristics

| Name | Scope | # concepts | Median | Subs. Hier | Version |
|------|-------|-----------|--------|-----------|---------|
| SNOMED CT | Clinical medicine (patient records) | 310,314 | 2 | yes | July 31, 2007 |
| LOINC | Clinical observations and laboratory tests | 46,406 | 3 | no | Version 2.21 (no "natural language" names) |
| FMA | Human anatomical structures | ~72,000 | ? | yes | (not yet in the UMLS) |
| Gene Ontology | Functional annotation of gene products | 22,546 | 1 | yes | Jan. 2, 2007 |
| RxNorm | Standard names for prescription drugs | 93,426 | 1 | no | Aug. 31, 2007 |
| NCI Thesaurus | Cancer research, clinical care, public information | 58,868 | 2 | yes | 2007_05E |
| ICD-10 | Diseases and conditions (health statistics) | 12,318 | 1 | no | 1998 (tabular) |
| MeSH | Biomedicine (descriptors for indexing the literature) | 24,767 | 5 | no | Aug. 27, 2007 |
| UMLS . | Terminology integration in the life sciences | 1,4 M | 2 | n/a | 2007AC (English only) |

NLM

[Bodenreider, YBMI 2008]

# NCI Thesaurus

# NCI thesaurus Characteristics (1)

◆ Current version: 08.08d (~monthly releases)

◆ Type: Controlled terminology / ontology

◆ Domain: Cancer

◆ Developer: NCI Center for Bioinformatics

◆ Funding: NCI

◆ Availability

  ● Publicly available: Yes

  ● Repositories: UMLS / OBO / NCBO BioPortal

◆ URL: http://nciterms.nci.nih.gov/

# NCI thesaurus Characteristics (2)

◆ Number of
- Concepts: 58,868 (2007_05E)
- Terms: 2.68 per concept

◆ Major organizing principles:
- Subsumption hierarchy
- Rich set of associative relationships
- Small proportion of defined concepts (many primitives)
- Links to many external resources

◆ Formalism: OWL Lite

# NCI thesaurus  Top level

## NCI_Thesaurus Taxonomy

- Abnormal Cell
- Activity
- Anatomic Structure, System, or Substance
- Biochemical Pathway
- Biological Process
- Chemotherapy Regimen or Agent Combination
- Conceptual Entity
- Diagnostic, Therapeutic, and Research Equipment
- Diagnostic or Prognostic Factor
- Disease, Disorder or Finding
- Drug, Food, Chemical or Biomedical Material
- Experimental Organism Anatomical Concept
- Experimental Organism Diagnosis
- Gene
- Gene Product
- Molecular Abnormality
- NCI Administrative Concept
- Organism
- Property or Attribute
- Retired Concept

# NCI thesaurus Example

# *Unified Medical Language System (UMLS)*

# UMLS Characteristics (1)

◆ Current version: 2008AA (2-3 annual releases)

◆ Type: Terminology integration system

◆ Domain: Biomedicine

◆ Developer: NLM

◆ Funding: NLM (intramural)

◆ Availability

 ● Publicly available: Yes* (cost-free license required)

 ● Repositories: UMLS

◆ URL: http://umlsks.nlm.nih.gov/

# UMLS Characteristics (2)

- Number of
  - Concepts: 1.5M (2008AA)
  - Terms: ~6M
- Major organizing principles (Metathesaurus):
  - Concept orientation
  - Source transparency
  - Multi-lingual through translation
- Formalism: Proprietary format (RRF)

# Addison's Disease: Concept

Disease or Syndrome

Addison's Disease

SNOMED CT
SNOMED Intl
MeSH
MedDRA
...

C0001403

ADRENAL INSUFFICIENCY (ADDISON'S DISEASE)
ADRENOCORTICAL INSUFFICIENCY, PRIMARY FAILURE
Hypoadrenalisms, Primary
Melasma addisonii
Primary adrenal deficiency
Asthenia pigmentosa
Bronzed disease
Insufficiency, adrenal primary
Primary adrenocortical insufficiency
Addison's, disease

Maladie d'Addison - French
Addison-Krankheit - German
Morbo di Addison - Italian
Doença de Addison - Portuguese
АДДИСОНОВА БОЛЕЗНЬ - Russian
アジソン病 - Japanese

An adrenal disease characterized by the progressive destruction of the adrenal cortex, resulting in insufficient production of aldosterone and hydrocortisone. Clinical symptoms include anorexia; nausea; weight loss; muscle ewakness; and hyperpigmentation of the skin due to increase in circulating levels of ACTH precursor hormone which stimulates melanocytes.

# Biomedical ontologies in action

[Bodenreider, YBMI 2008]

- ◆ Functional perspective
  - ● What are they for (vs. what are they)?
- ◆ "High-impact" biomedical ontologies
- ◆ 3 major categories of use
  - ● Knowledge management (indexing and retrieval of data and information, access to information, mapping among ontologies)
  - ● Data integration, exchange and semantic interoperability
  - ● Decision support and reasoning (data selection and aggregation, decision support, natural language processing applications, knowledge discovery).

# Ontologies for mining biomedical data

*Normalization*
*Integration*
*Aggregation*

# Ontologies for mining biomedical data

*Normalization*
*Integration*
*Aggregation*

# Issues

◆ Variability of natural language

- Lexical variants
  - Lung cancer
  - Cancer of the lung
  - Lung cancers

- Synonyms
  - Pulmonary cancer
  - Malignant neoplasm of lung
  - Malignant tumor of lung

# Solutions with ontologies

◆ Controlled vocabulary

  ● Standard list of terms to be used for a given purpose

    ▪ Gene Ontology (functional annotation of gene products)

    ▪ MeSH (indexing of biomedical articles)

    ▪ ICD (mortality and morbidity reporting)

# London Bills of Mortality

# Solutions

- ◆ Controlled vocabulary
  - ● Standard list of terms to be used for a given purpose
    - ▪ Gene Ontology (functional annotation of gene products)
    - ▪ MeSH (indexing of biomedical articles)
    - ▪ ICD (mortality and morbidity reporting)
- ◆ Lexical normalization programs [McCray, SCAMC 1994]
  - ● Management of lexical terminological variability
  - ● UMLS *Lexical Variant Generation* program
  - ● Used in terminology integration systems (e.g., UMLS)
  - ● Useful for indexing and text mining applications

# Normalization

| Process | | Result |
|---|---|---|
| | | Hodgkin's diseases, NOS |
| Remove genitive | → | Hodgkin diseases, NOS |
| Remove stop words | → | Hodgkin diseases, |
| Lowercase | → | hodgkin diseases, |
| Strip punctuation | → | hodgkin diseases |
| Uninflect | → | hodgkin disease |
| Sort words | → | disease hodgkin |

NLM

25

# Normalization: Example

Hodgkin Disease
HODGKINS DISEASE
Hodgkin's Disease
Disease, Hodgkin's
Hodgkin's, disease
HODGKIN'S DISEASE
Hodgkin's disease
Hodgkins Disease
Hodgkin's disease NOS
Hodgkin's disease, NOS
Disease, Hodgkins
Diseases, Hodgkins
Hodgkins Diseases
Hodgkins disease
hodgkin's disease
Disease, Hodgkin

normalize ⟶ disease hodgkin

# Ontologies for mining biomedical data

*Normalization*

*Integration*

*Aggregation*

# Issues

◆ Different codes for the same biomedical entity in different ontologies

- SNOMED CT:    363732003    Addison's disease
- MeSH:             D000224        Addison Disease
- NCI Thesaurus:  C26689        Addison's Disease
- ICD 9-CM:         255.41          Addison's disease NOS
- ICD 10:            E27.1            Primary adrenocortical insufficiency
- MedDRA:          10001130     Addison's disease
- …

◆ Hindrance to the integration of datasets (e.g., clinical, research and epidemiology data)

# Solutions with ontologies

◆ Identify equivalent concepts across ontologies

◆ Specific mappings
  ● SNOMED to ICD 9-CM (provided by SNOMED)

◆ Terminology integration systems
  ● Manually curated
    ▪ Unified Medical Language System (UMLS)
    ▪ RxNorm (for drug vocabularies)
  ● Automatically aligned
    ▪ BioPortal

# Terminology integration in the UMLS

- SNOMED CT:    363732003    Addison's disease
- MeSH:    D000224    Addison Disease
- NCI Thesaurus:  C26689    Addison's Disease
- ICD 9-CM:    255.41    Addison's disease NOS
- ICD 10:    E27.1    Primary adrenocortical insufficiency
- MedDRA:    10001130    Addison's disease
- …

C0001403

◆ Identified as synonyms (semi-automatically)

◆ Clustered into a UMLS concept

◆ Assigned a permanent identifier (CUI)

# Integrating subdomains



31

# Integrating subdomains



Clinical repositories

Genetic knowledge bases

Other subdomains

Biomedical literature

Model organisms

Genome annotations

Anatomy

NLM

32

# Trans-namespace integration

# Ontologies for mining biomedical data

*Normalization*

*Integration*

*Aggregation*

# Issues

◆ Various levels of granularity
  - Upper limb
    - Hand
      – Index finger
        » Diaphysis of distal phalanx of left index finger

◆ Fine-grained  may not be appropriate for high-level analysis
  - Reduce statistical power

◆ Need to abstract away from details
  - Aggregate into a more generic concept
  - Corollary: Apply to more specific concepts

35

# Solutions with ontologies

◆ Aggregate along subsumption hierarchies
◆ Helps enrich feature sets for data mining purposes

◆ Examples
   ● GO Slims (analysis of functional annotations)
   ● Categorization of adverse events based on high-level disease categories
◆ Corollary
   ● Patient selection based on high-level ICD 9-CM codes
   ● MeSH term "explosion" (information retrieval)

# Aggregation with MeSH

# Ontologies for mining biomedical data

*Text mining*
*Analysis of gene expression data*
*Mining adverse drug reactions*

# Ontologies for mining biomedical data

*Text mining*

*Analysis of gene expression data*
*Mining adverse drug reactions*

# Ontological resources for text mining

◆ Lexical resources
  - SPECIALIST lexicon (UMLS) / LVG
  - Lexico-syntactic analysis, normalization
◆ Terminological resources
  - UMLS Metathesaurus / MetaMap
  - Named entity recognition, semantic normalization
◆ Ontological resources
  - UMLS Semantic Network / SemRep
  - Relation extraction, semantic interpretation

[Ananiadou, Text mining for biology and biomedicine 2006]

40

# Ontologies for mining biomedical data

*Text mining*

*Analysis of gene expression data*

*Mining adverse drug reactions*

# Traditional approach

Analysis of gene expression data

◆ Cluster analysis
 - Genes
 - Genes and samples

◆ Elicitation of clusters using external knowledge
 - Functional annotations
 - Participation in pathways

42

# Clustering constraints from ontologies

◆ Use ontologies as a source of prior knowledge

◆ Ontologies used to constrain the clustering process

  ● Several variants of the clustering algorithms

◆ Tends to result in more meaningful clusters

[Liu, CSB 2004]
[Kustra, CBMS 2006]
[Huang, Omics 2006]
[Chabalier, BMC Bioinfo 2007]

43

# Ontologies for mining biomedical data

*Text mining*
*Analysis of gene expression data*
*Mining adverse drug reactions*

# Identifying adverse drug reactions

- ◆ Pharmacovigilance of self-reported ADR cases
  - Coded with MedDRA
  - Manually curated
- ◆ Bayesian analysis of the drug-ADR associations
- ◆ 4 variants
  - MedDRA without subsumption links
  - MedDRA with original subsumption links
  - MedDRA with enhanced subsumption links
  - MedDRA with enhanced subsumption links and approximate matching [Henegar, CBM 2006]

45

# Using subsumption links increases the signal



[Henegar, CBM 2006]

Enhanced aggregation

Basic aggregation

- Ontology TR & AM
- Ontology TR
- MedDRA TR
- MedDRA without TR

# Conclusions

# Translational research  NIH Roadmap

# Clinical and Translational Science Awards

49

# Ontologies for data mining

◆ Ontologies

- Normalize datasets
- Aggregate data of different granularity

increase signal

◆ Ontology integration systems

- Integrate datasets

◆ Challenges

- Permanent identifiers for biomedical entities
- Availability
- Quality

# Data mining with ontologies

- ◆ Ontologies are increasingly used in biological data mining
  - Text mining
    - Named entity recognition
    - Relation extraction
  - In combination with other features
  - To enhance feature sets

- ◆ Few data mining algorithms natively take advantage of ontologies

# Medical Ontology Research

Contact: olivier@nlm.nih.gov
Web: mor.nlm.nih.gov



*Olivier Bodenreider*

Lister Hill National Center
for Biomedical Communications
Bethesda, Maryland - USA