Mini Course From Data Science to interpretable Al





Assoc. Prof. Dr. Andreas HOLZINGER (Medical University Graz)

Day 1 > Part 1 > Monday, 17.06.2019

Data, Information and Knowledge Representation

This is the version for printing and reading. The lecture version is didactically different.

Health Informatics - Andreas Holzinger

-

Overview



Day 1 - Fundamentals

01 Introduction to Al Machine Learning

02 Data, Information & Knowledge Representation

03 Decision Making and Decision Support

04 Causal Reasoning and Interpretable Al

Keywords



- Data
- Information
- Knowledge
- Dimensionality of data
- Biomedical Ontologies
- Standardized Medical Data
- SNOMED
- UMLS



- ... be aware of the types and categories of different data sets in biomedical informatics;
- ... know some differences between data, information, and knowledge;
- ... be aware of standardized/non-standardized and well-structured/"un-structured" information/data;
- ... have a basic overview on some ontological approaches for standardized medicine;
- ... have some background on classifications

5

Advance Organizer (2/2)



- Induction = deriving a likely general conclusion from a set of particular statements;
- Information = derived from the data by interpretation (with feedback to the clinician);
- Information Entropy = a measure for uncertainty: highly structured data contain low entropy, if everything is in order there is no uncertainty, no surprise, ideally H = 0
- Knowledge = obtained by inductive reasoning with previously interpreted data, collected from many similar patients or processes, which is added to the "body of knowledge" (explicit knowledge). This knowledge is used for the interpretation of other data and to gain implicit knowledge which guides the clinician in taking further action;
- Large Data = consist of at least hundreds of thousands of data points
- Multi-Dimensionality = containing more than three dimensions and data are multivariate
- Multi-Modality = a combination of data from different sources
- Multivariate = encompassing the simultaneous observation and analysis of more than one statistical variable;
- Reasoning = process by which clinicians reach a conclusion after thinking on all facts;
- Spatiality = contains at least one (non-scalar) spatial component and non-spatial data
- Structural Complexity = ranging from low-structured (simple data structure, but many instances, e.g., flow data, volume data) to high-structured data (complex data structure, but only a few instances, e.g., business data)
- Time-Dependency = data is given at several points in time (time series data)
- Voxel = volumetric pixel = volumetric picture element

Advance Organizer (1/2)



- Abduction = cyclical process of generating possible explanations (i.e., identification of a set of hypotheses that are able to account for the clinical case on the basis of the available data) and testing those (i.e., evaluation of each generated hypothesis on the basis of its expected consequences) for the abnormal state of the patient at hand;
- Abstraction = data are <u>filtered according to their relevance</u> for the problem solution and chunked in schemas representing an abstract description of the problem (e.g., abstracting that an adult male with haemoglobin concentration less than 14g/dL is an anaemic patient):
- Artefact/surrogate = <u>error</u> or <u>anomaly</u> in the perception or representation of information trough the involved method, equipment or process;
- Data = physical entities at the lowest abstraction level which are, e.g. generated by a
 patient (patient data) or a (biological) process; data contain no meaning;
- Data quality = Includes quality parameter such as: Accuracy, Completeness, Update status, Relevance, Consistency, Reliability, Accessibility;
- Data structure = way of storing and <u>organizing</u> data to use it <u>efficiently;</u>
- Deduction = deriving a particular valid conclusion from a set of general premises;
- DIK-Model = Data-Information-Knowledge three level model
- Disparity = containing different types of information in different dimensions
- Heart rate variability (HRV) = measured by the variation in the beat-to-beat interval;
- HRV artifact = noise through errors in the location of the instantaneous heart beat, resulting in errors in the calculation of the HRV, which is highly sensitive to artifact and errors in as low as 2% of the data will result in unwanted biases in HRV calculations;

Health Informatics - Andreas Holzinger

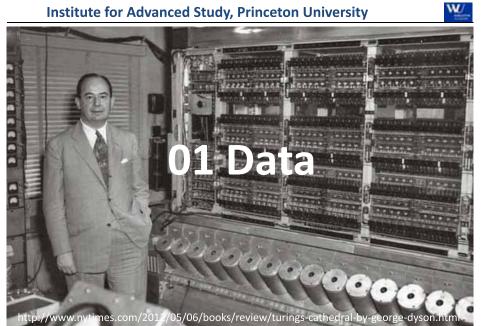
6

Agenda



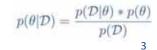
- 00 Reflection follow-up from last lecture
- 01 What is data?
- 02 On Standardization
- 03 Knowledge Representation
- 04 Biomedical Ontologies
- 05 Medical Classifications

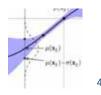




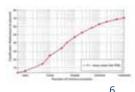


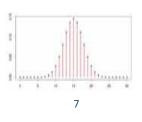


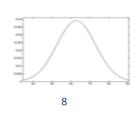


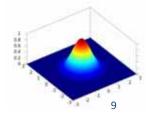












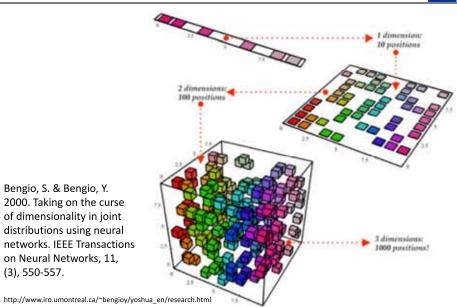
Traditional Statistics versus Machine Learning



- Data in traditional **Statistics**
- Low-dimensional data ($< \mathbb{R}^{100}$)
- Problem: Much noise in the data
- Not much structure in the data but it can be represented by a simple model

- Data in Machine Learning
- High-dimensional data ($\gg \mathbb{R}^{100}$)
- Problem: not noise, but complexity
- Much structure, but the structure can not be represented by a simple model

Health Informatics - Andreas Holzinger



Holginger, A., Dehmer, M. & Jurisica, 1. 2014. Knowledge Discovery and interactive Data Mining in Bioinformatics - State-of-the-Art, future challenges and research directions. BMC Bioinformatics, 15, (56), II.

Collective

Individual

Tissue

Cell

Bacteria

Virus

Molecule

Atom

Health Informatics - Andreas Holzinger

1

Data for clinical purposes – integration is unsolved!



Private Health vault data Electronic health record data Physiological data Laboratory results

Metabolomics Chemical processes Cellular reactions Enzymatic reactions

Metabolomics Chemical processes Cellular reactions Enzymatic reactions

Proteomics
Protein-Protein Interactions

Epigenetics
Epigenetic modifications

Exposome Environmental data Air pollution Exposure (toxicants)



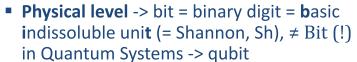
Collective data
Social data
Fitness, Wellness data
Ambient Assisted Living data
(Non-medical) personal data

Foodomics, Lipidomics Nutrition data (Nutrigenomics) Diet data (allergenics)

Imaging data X-Ray, ultrasound, MR, CT, PET, cams, observation (e.g. sleep laboratory), gait (child walking)

Transcriptomics RNA, mRNA, rRNA, tRNA

Taxonomy of data



- Logical Level -> integers, booleans, characters, floating-point numbers, alphanumeric strings, ...
- Conceptual (Abstract) Level -> data-structures, e.g. lists, arrays, trees, graphs, ...
- Technical Level -> Application data, e.g. text, graphics, images, audio, video, multimedia, ...
- "Hospital Level" -> Narrative (textual) data, genetic data, numerical measurements (physiological data, lab results, vital signs, ...), recorded signals (ECG, EEG, ...), Images (cams, x-ray, MR, CT, PET, ...)



Examples: Imaging Data



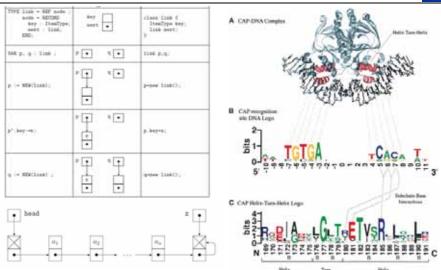


Health Informatics - Andreas Holzinger

17

Example Data Structures (1/3): List





Crooks, G. E., Hon, G., Chandonia, J. M. & Brenner, S. E. (2004) WebLogo: A sequence logo generator. Genome Research, 14, 6, 1188-1190.

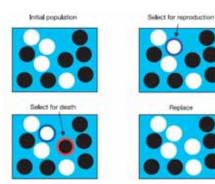
Health Informatics - Andreas Holzinger

15

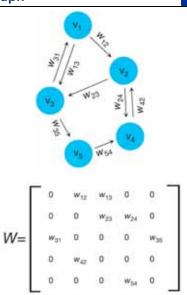
Example Data Structures (2/3): Graph



Evolutionary dynamics act on populations. Neither genes, nor cells, nor individuals evolve; only populations evolve.



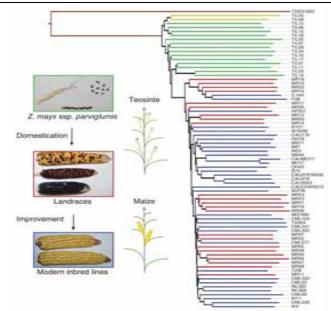
Lieberman, E., Hauert, C. & Nowak, M. A. (2005) Evolutionary dynamics on graphs. *Nature*, *433*, *7023*, *312-316*.



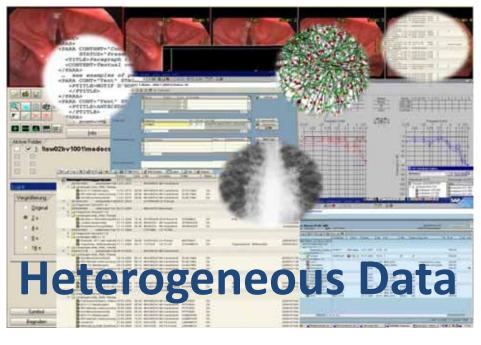
Example Data Structures (3/3) Tree



Hufford et. al. 2012. Comparative population genomics of maize domestication and improvement. *Nature Genetics*, 44, (7), 808-811.







21

Natural Language is a good example for complexity ...



Problem: Context!



Biomedical R&D data (e.g. clinical trial data) Clinical patient data (e.g. EPR, lab, reports etc.)

The combining link is text

Health business data (e.g. costs, utilization, etc.

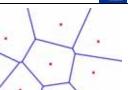
Private patient data (e.g. AAL, monitoring, etc.)

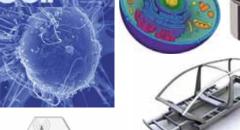
Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C. & Byers, A. H. (2011) *Big data: The next frontier for innovation, competition, and productivity. Washington (DC), McKinsey Global Institute.*

Health Informatics - Andreas Holzinger

22

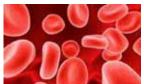
Semantic Ambiguity – Missing Context











Health Informatics – Andreas Holzinger 23 Health Informatics – Andreas Holzinger



Is a picture really worth a thousand words?

Health Informatics - Andreas Holzinger

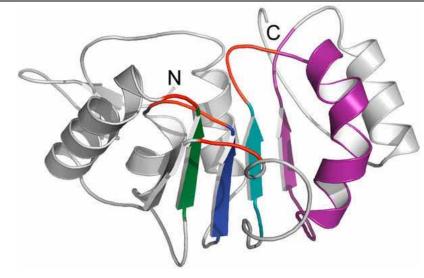
25







Holzinger, A., Geierhofer, R. & Errath, M. 2007. Semantische Informationsextraktion in medizinischen Informationssystemen. *Informatik Spektrum*, *30*, *(2)*, *69-78*.



Magnani, R., et al. 2010. Calmodulin methyltransferase is an evolutionarily conserved enzyme that trimethylates Lys-115 in calmodulin. *Nature Communications*, *1*, *43*.

Health Informatics - Andreas Holzinger

26



02 Medical Communication

Health Informatics – Andreas Holzinger

27

Health Informatics - Andreas Holzinger

28



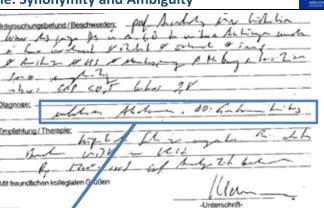


Holzinger, A., Geierhofer, R., Ackerl, S. & Searle, G. (2005). *CARDIAC@VIEW: The User Centered Development of a new Medical Image Viewer. Central European Multimedia and Virtual Reality Conference, Prague, Czech Technical University (CTU), 63-68.*

Health Informatics - Andreas Holzinger

29

German Example: Synonymity and Ambiguity



"die Antrumschleimhaut ist durch Lymphozyten infiltriert" "lymphozytäre Infiltration der Antrummukosa" "Lymphoyteninfiltration der Magenschleimhaut im Antrumbereich"

	Anfo: NCHIN
Kurzanamnese: St.p. SHT	
Fragestellung: -	
	pecial Words
SB	
Bewegungsartefakte. Zustand nach Schädelhimtrauma.	anguage Mix
Das Cor in der Größennorm, keine akuten Stauungszeichen. Fragliches Infiltrat parahilär li. im UF, RW-Erguss li.	hbreviations
Fragliches Infiltrat parahillär li. im UF, RW-Erguss li. Zustand nach Anlage eines ET, die Spitze ca. 5cm cranial der Bir positioniert. ZVK über re., die Spitze in Proj. auf die VCS. Kein Hi Der re. Rezessus frei.	furkation, lieg. MS, orthotop inweis auf Pneumothorax. TTOTS
Mit kollegialen Grüßen	
*** Elektronische Freigabe durch am 09.0	05.2006

Holzinger, A., Geierhofer, R. & Errath, M. 2007. Semantische Informationsextraktion in medizinischen Informationssystemen. *Informatik Spektrum, 30, (2), 69-78*.

Health Informatics - Andreas Holzinger

30

German Local Hospital Abbreviations ... (example)



- HWI =
 - Harnwegsinfekt
 - Hinterwandinfarkt
 - Hinterwandischämie
 - Hakenwurminfektion
 - Halswirbelimmobilisation
 - Hip Waist Index
 - Height-Width Index
 - Heart-Work Index
 - Hemodynamically weighted imaging
 - High Water Intake
 - Hot water irrigation
 - Hepatitic weight index
 - Häufig wechselnder Intimpartner
- Leitung = Nervenleitung, Abteilungsleitung, Stromleitung, Wasserleitung, Harnleitung, Ableitung, Vereinsleitung ©...



- Syntax
- Semantics
- Pragmatics
- Context
- [Emotion]



Andrej Karpathy & Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 3128-3137.

33



Key Challenges

- Increasingly large data sets due to data-driven medicine [1]
- Increasing amounts of non-standardized data and un-structured information (e.g. "free text")
- Data quality, data integration, universal access
- Privacy, security, safety, data protection, data ownership, fair use of data [2]
- **Time** aspects in databases [3]

[1] Shah, N. H. & Tenenbaum, J. D. 2012. The coming age of data-driven medicine: translational bioinformatics' next frontier. Journal of the American Medical Informatics Association, 19, (E1), E2-E4. [2] Kieseberg, P., Hobel, H., Schrittwieser, S., Weippl, E. & Holzinger, A. 2014. Protecting Anonymity in Data-Driven Biomedical Science. In: LNCS 8401. Berlin Heidelberg: Springer pp. 301-316... [3] Gschwandtner, T., Gärtner, J., Aigner, W. & Miksch, S. 2012. A taxonomy of dirty time-oriented data. In: LNCS 7465. Heidelberg, Berlin: Springer, pp. 58-72.



Thomas, J. J. & Cook, K. A. 2005. Illuminating the path: The research and development agenda for visual analytics, New York, IEEE Computer Society Press.

SEMANTICS SYNTAX ORPHOLOGI PONOLOG HONETICA **Linguistic Data** rases and sente meaning of phrases and meaning in context of discov

PRAGMATICS

Health Informatics - Andreas Holzinger

Health Informatics - Andreas Holzinger



03 On **Standardization**



Still a big problem: Inaccuracy of medical data



- Medical (clinical) data are defined and detected disturbingly "soft" ...
- ... having an obvious degree of variability and inaccuracy.
- Taking a medical history, the performance of a physical examination, the interpretation of laboratory tests, even the definition of diseases ... are surprisingly **inexact**.
- Data is defined, collected, and interpreted with a degree of variability and inaccuracy which falls far short of the standards which engineers do expect from most data.
- Moreover, standards might be interpreted variably by different medical doctors, different hospitals, different medical schools, different medical cultures. ...

Komaroff, A. L. (1979) The variability and inaccuracy of medical data. Proceedings of the IEEE, 67, 9, 1196-1207.

Quest for standardization as old as med, informatics



IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING, VOL. BME-19, NO. 5, SEPTEMBER 1972

HEWLETT-PACKARD LIBRARY331

Standardization and Health Care AUG 18 1972

J. H. U. BROWN, SENIOR MEMBER, IEEE, AND DEWITT JAMES LOWELD. Not Reserve

Abstract-In order to deliver reasonable health care to all people, it is essential that standards be established. Standards vary with the type of control and with the approach desired in determining the quality of care. This paper discusses various kinds of standards and their application in the health care field. Standards may be determined as a cess or as a direct regulation. It is probable that regulation of stanfards by process is the most satisfactory method.

arbiter may be the market place or agencies that rely on expertise from many sources to set acceptable standards of quality or performance. For these reasons, the final moderator may be found in a governmental authority, and its delegation into a system of regulation, law, and judicial action, so that an established code can become the focal point of resolution.

INTRODUCTION

COCIETY cannot exist without a yardstick by which its access and output by virtue of the goal and process structure

THE OBJECTIVES OF STANDARDIZATION

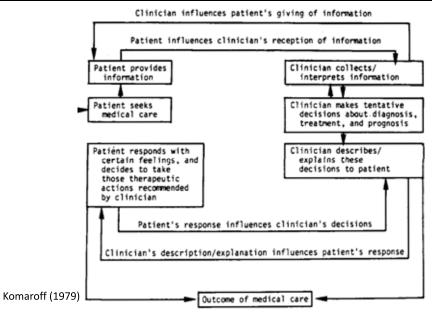
complishments or failures are measured. Such yardsticks tablish quality. However, they accomplish more for society are called standards. They are created by the need for regulation and control as an escape from anarchy or to motivate mance. A standard allows coordination of effort between towards greater achievement. In the ultimate, society dictates producers so that like products can be produced. It permits these limits by the demands it places upon itself. Standards the reproduction of similar units in mass quantity and permits provide opportunities for security and augmentation of proby performance. It establishes freedom of interchange of material and ideas, and permits the activity in one part of society

Brown, J. H. U. & Loweli, D. J. (1972) Standardization and Health Care. IEEE Transactions on Biomedical Engineering, BME-19, 5, 331-334.

Health Informatics - Andreas Holzinger

The patient-clinician dialogue (from 1979)





Health Informatics - Andreas Holzinger Health Informatics - Andreas Holzinger



- ... ensures that information is interpreted by all users with the same understanding;
 - supports the reusability of the data,
 - improves the efficiency of healthcare services and
 - avoids errors by reducing duplicated efforts in data entry;
- Data standardization refers to
 - a) the data content;
 - b) the terminologies that are used to represent the data;
 - c) how data is exchanged; and
 - iv) how knowledge, e.g. clinical guidelines, protocols, decision support rules, checklists, standard operating procedures are represented in the health information system (refer to IOM).
- Elements for sharing require standardization of identification, record structure, terminology, messaging, privacy etc.
- The most used standardized data set to date is the International Classification of Diseases (ICD), which was first adopted in 1900 for collecting statistics (Ahmadian et al. 2011)

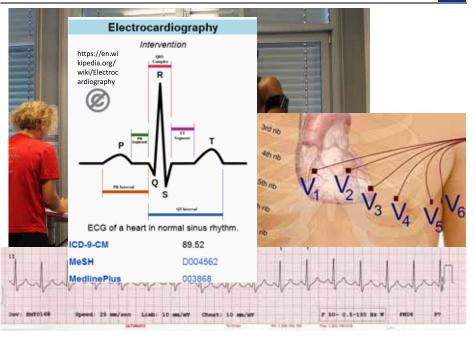
41

Standardization of ECG data (1/2)



- There has been a large number of ECG storage formats proclaiming to promote interoperability.
- There are three predominant ECG formats:
 - SCP-ECG (1993, European Standard, Binary data)
 - DICOM-ECG (2000, European Standard, Binary data)
 - HL7 aECG (2001, ANSI Standard, XML data)
- A mass of researchers have been proposing their own ECG storage formats to be considered for implementation (= proprietary formats).
- Binary has been the predominant method for storing ECG data

Bond, R. R., Finlay, D. D., Nugent, C. D. & Moore, G. (2011) A review of ECG storage formats. *International Journal of Medical Informatics*, 80, 10, 681-697.



Standardization of ECG (2/2)

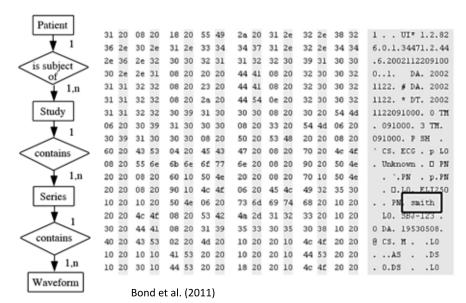


Overview on current ECG storage formats

ECG format	Year	Method of implemen- tation	Specification	Viewers
SCP-ECG	1993	BINARY	Can be freely downloaded from the Internet [7].	Freely available SCP-ECG Viewer made by EcgSoft [8]
DICOM-WS 30	2000	BINARY	Can be freely downloaded from the Internet [5].	Freely available DICOM-ECG viewer made by Charruasoft [9].
HL7 aECG	2001	XML	The XML Schema can be used as the specification or the implementation guide by AMPS [6].	Freely available aECG viewer by AMPS [10].
ecgML	2003	XML	Can be freely downloaded from the Internet [11].	None currently exist. Under development.
MFER	2003	BINARY	Can be freely downloaded from the Internet [12].	Freely available MFER viewer [13].
Philips XML	2004	XML	The specification is packaged with the actual product.	Philips viewer. Not freely available.
XML-EOG	2007	XML	Can be freely downloaded from the Internet [14].	XML-ECG viewer [14]. Not freely available.
mitCGml	2008	XML	Can be freely downloaded from the Internet [15].	mECGml mobile viewer [15]. Not freely available.
ecgAware	2008	XML	Can be freely downloaded from the Internet [16].	TeleCardio viewer [16]. Not freely available.

Bond, R. R., Finlay, D. D., Nugent, C. D. & Moore, G. (2011) A review of ECG storage formats. *International Journal of Medical Informatics*, 80, 10, 681-697.





4



04 Knowledge Representation



Bond et al. (2011)

Health Informatics - Andreas Holzinger

Health Informatics - Andreas Holzinger

Examples for famous knowledge representations



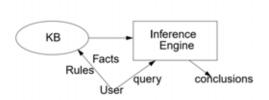
Mathematical Logic	Psychology	Biology	Statistics	Economics
Aristotle	100			
Descartes				
Boole	James		Laplace	Bentham Pareto
Frege			Bernoullii	Friedman
Peano				
and the same of th	Hebb	Lashley	Bayes	
Goedel	Bruner	Rosenblatt		
Post	Miller	Ashby	Tversky,	Von Neumann
Church	Newell,	Lettvin	Kahneman	Simon
Turing	Simon	McCulloch, Pitts		Raiffa
Davis		Heubel, Weisel		
Putnam				
Robinson				
Logic soas		Connectionism	Causal	Rational
PROLOG KBS	Frames		Networks	Agents

Davis, R., Shrobe, H., Szolovits, P. 1993 What is a knowledge representation? Al Magazine, 14, 1, 17-33.



Formalization versus Expressivity







Logical reasoning can be dangerous:

A dime is better than a nickel.
 A nickel is better than a penny.
 Therefore, a dime is better than a penny.
 A penny is better than a nothing.
 Nothing is better than world peace.
 Therefore, a dime is better than a penny.

Health Informatics - Andreas Holzinger

49

Example for Modeling of biomedical knowledge





Hajdukiewicz, J. R., Vicente, K. J., Doyle, D. J., Milgram, P. & Burns, C. M. (2001) Modeling a medical environment: an ontology for integrated medical informatics design. *International Journal of Medical Informatics*, 62, 1, 79-99.

Blobel, B. (2011) Ontology driven health information	Formal ontologies	Di Propi	General logic Modal logic First-order logic Description logic spositional logic Il languages	
	Meta-data and data models	Formal taxono Data models XML Schema Database schemas	mies	
systems architectures enable pHealth for empowered		Principled, informational h ML DTD ctured glossaries ri	Thesauri and taxonomies	
patients. International Journal of Medical	Data dictio Ad hoc hierarc "ordinary" glossal Terms	hies	Glossaries and data dictionaries	
Informatics, 80, 2, e17-e25.	-		Formalization	

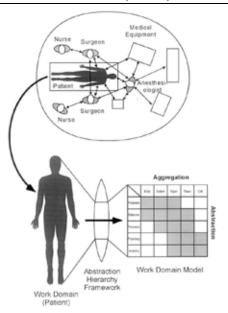
Health Informatics - Andreas Holzinger

50

Building and Creating a work domain model (WDM)







Mass

Inflow

Blood

Infusion

Adequate Blood Volume

Mass

Store

Heart Function

(HR, rhythm,

contractility)

Pulmonary

Circulation

53

Cardiovascular Example

Mass

Transfer

Systemic

Circulation

Adequate Circulation

Mass

Outflow

Blood

Volume

Large vein

Blood

Loss

Artery

WHY?

WHAT?

HOW?

Level of

Abstraction

Purposes

Balances

Processes

Physiology

Anatomy

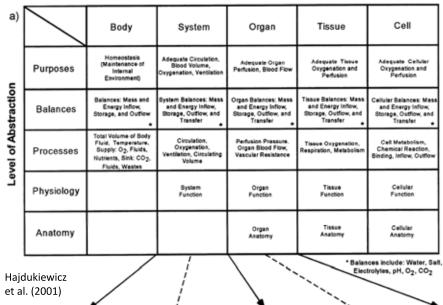
Haidukiewicz

Health Informatics - Andreas Holzinger

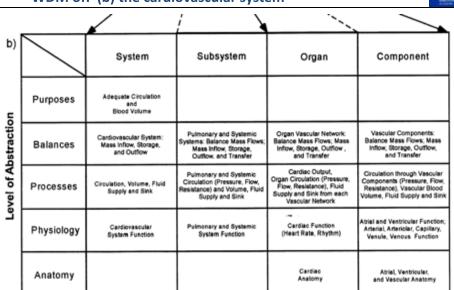
et al. (2001)



Level of Aggregation



WDM of: (b) the cardiovascular system

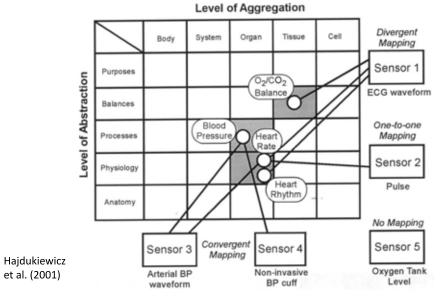


Hajdukiewicz et al. (2001)

Example: Mapping OR sensors onto the WDM

et al. (2001)

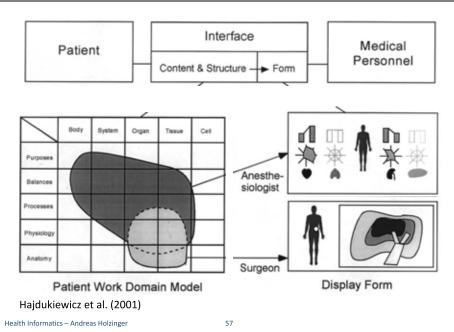








W/



05 Ontologies

Health Informatics - Andreas Holzinger

5

A simple question: What is a Jaguar?









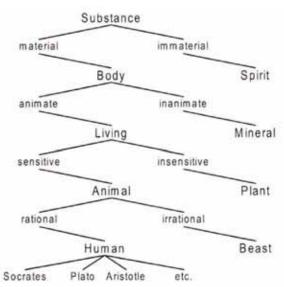


The first "Ontology of what exists"



* 384 BC † 322 BC

Simonet, M., Messai, R., Diallo, G. & Simonet, A. (2009) Ontologies in the Health Field. In: Berka, P., Rauch, J. & Zighed, D. A. (Eds.) Data Mining and Medical Knowledge Management: Cases and Applications. New York, Medical Information Science Reference, 37-56.



Health Informatics – Andreas Holzinger 59 Health Informatics – Andreas Holzinger



- **Example: Network-Extracted Ontology of human cell**
- W

- Aristotle attempted to classify the things in the world where it is employed to describe the existence of beings in the world;
- Artificial Intelligence and Knowledge Engineering deals also with reasoning about models of the world.
- Therefore, AI researchers adopted the term 'ontology' to describe what can be computationally represented of the world within a program.
- "An ontology is a formal, explicit specification of a shared conceptualization".
 - A 'conceptualization' refers to an abstract model of some phenomenon in the world by having identified the relevant concepts of that phenomenon.
 - 'Explicit' means that the type of concepts used, and the constraints on their use are explicitly defined.

Studer, R., Benjamins, V. R. & Fensel, D. (1998) Knowledge Engineering: Principles and methods. *Data & Knowledge Engineering*, 25, 1-2, 161-197.

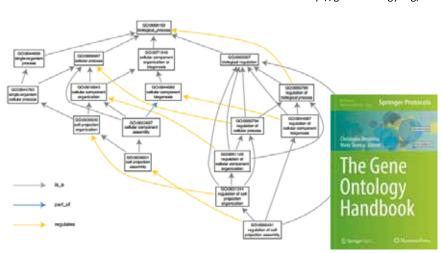
Health Informatics - Andreas Holzinger

61

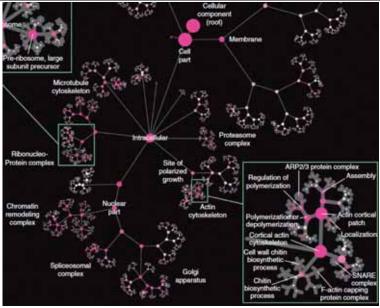
W/

Example: GO

http://geneontology.org/



Hastings, J. 2017. Primer on Ontologies. In: Dessimoz, C. & Škunca, N. (eds.) The Gene Ontology Handbook. New York, NY: Springer New York, pp. 3-13, doi:10.1007/978-1-4939-3743-1 1.



http://www.kurzweilai.net/images/cell-model.png (Credit: UC San Diego School of Medicine)

Ontology: Terminology



- Ontology = a structured description of a domain in form of concepts ↔ relations;
- The IS-A relation provides a taxonomic skeleton;
- Other relations reflect the domain semantics;
- Formalizes the terminology in the domain;
- Terminology = terms definition and usage in the specific context;
- Knowledge base = instance classification and concept classification;
- Classification provides the domain terminology

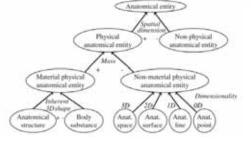
...

Additionally an ontology may satisfy:



- (1) In addition to the IS-A relationship, partitive (meronomic) relationships may hold between concepts, denoted by PART-OF. Every PART-OF relationship is irreflexive, asymmetric and transitive. IS-A and PART-OF are also called hierarchical relationships.
- (2) In addition to hierarchical relationships, associative relationships may hold between concepts. Some associative relationships are domain-specific (e.g., the branching relationship between arteries in anatomy and rivers in geography).
- (3) Relationships r and r' are inverses if, for every pair of concepts x and y, the relations (x, r, y) and (y, r', x) hold simultaneously. A symmetric relationship is its own inverse. Inverses of hierarchical relationships are called INVERSE-IS-A and HAS-PART, respectively.
- (4) Every non-taxonomic relation of x to z, (x, r, z), is either inherited ((y, r, z)) or refined ((y, r, z') where z' is more specific than z) by every child y of x. In other words, every child y of x has the same properties (z) as it parent or more specific properties (z').

Zhang, S. & Bodenreider, O. 2006. Law and order: Assessing and enforcing compliance with ontological modeling principles in the Foundational Model of Anatomy. *Computers in Biology and Medicine*, 36, (7-8), 674-693.



Health Informatics - Andreas Holzinger

6

Examples of Biomedical Ontologies

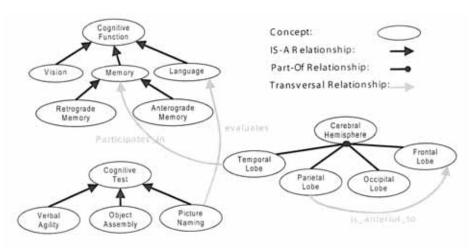


Nome Ref.	104	Scope		# concept names				Subs.	Version / Notes
	xope	concepts	Min	Max	Med	Avg	Hier.	Version / Notes	
SNOWED CT	[21]	Onical medicine (patient records)	310,314	n	37	2	2.57	yes	July 31, 2007
LOUNC	[24]	Genical observations and lichwartery tests	46,406	1	3	3	2.85	ns.	Version 2.21 (no "outural language" names)
EWA	[25]	Numer creaturaled structures	~72,000	10	2:	020	-1.50	yes	(not yet in the UMLS)
Gene Ontology	[28]	Functional annatation of gene products	22,546	-5	74	19	2.15	781	ton 3,7007
ReNorm	[31]	Standard names for prescription drugs	93,426	-1	2	- 1	1.10	30	Aug. 31, 2007
NCI Thesavas	[34]	Cancer research, clinical core, public information	58,868	1	100	2	2.68	76	2007_05E
KD-10	[26]	Disease and conditions (health statistics)	12,318	-17	31	- 12	1.00	86	1998 (habulus)
MeSH	[38]	Burnedicine (descriptors for indexing the literature)	24,767	Ŧ	208	5	7,47	14	Aug. 27, 2007
UMLS Mins.	[41]	Terrinology integration in the Me sciences	1,4.8	1	339	2	3.77	1/0	2007AC (English only)

Bodenreider, O. (2008) Biomedical ontologies in action: role in knowledge management, data integration and decision support. *Methods of Information In Medicine*, 47, Supplement 1, 67-79.

Example of a conceptual structure from CogSci





Simonet, M., Messai, R., Diallo, G. & Simonet, A. (2009) Ontologies in the Health Field. In: Berka, P., Rauch, J. & Zighed, D. A. (Eds.) *Data Mining and Medical Knowledge Management: Cases and Applications. New York, Medical Information Science Reference, 37-56.*

Health Informatics - Andreas Holzinger

66

Taxonomy of Ontology Languages



■ 1) Graph notations

- Semantic networks
- Topic Maps (ISO/IEC 13250)
- Unified Modeling Language (UML)
- Resource Description Framework (RDF)

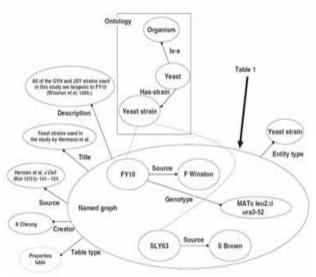
2) Logic based

- Description Logics (e.g., OIL, DAML+OIL, OWL)
- Rules (e.g. RuleML, LP/Prolog)
- First Order Logic (KIF Knowledge Interchange Format)
- Conceptual graphs
- (Syntactically) higher order logics (e.g. LBase)
- Non-classical logics (e.g. Flogic, Non-Mon, modalities)

3) Probabilistic/fuzzy

Example for (1) Graphical Notation: RDF





Name	Genotype*	Source
PYNE	MMT+ lmJA1 umJ-52	F Wasston
PY22	MAT's his IA200 unit FS2	F Wissten
CHYL	MMT's lim251 his 65,200 una 6-52 malm20-1	This study
HYTH?	MAT's hind \$200 and 52 gam lift (NB3)	This study
357948	MMPs len231/len231 until-52/until-52	This study
153/999	MAT's len2/A1 his/A200 una F-52	This study
3571065	MAT's les251 hts35200 unu3-52 mdm200:: LBIO	This study
3571094	MMPs length / his H\$200 unit F-52 gree ED: MEST	This study
3571138	MAT's les23-1-fes23-1 his/13/200 fis/35/200 unal-52/unal-52 special-8853/ + endex200-1-fes2/ +	This study
1911285	M647's len2/62 his/35/200 anal-52 (pm27): A85.1	This study
3571340	M64P+ lm252 his/35200 unu3-52 mdm20D1: LBSO	This study
380374	M64°s les231/fes231 fes33200 fes33200 ses23 52'ses23 52 spes201 8853/ 4 mdm2001: LESU7 4	This study
ABT1249	MAE's load 3, 112 analy 52 lyad 401 adad 301 adad bersal 40	A Bretscher
BCY4	M64Ps loss2-3, 112 http://dx.200i.oru/1-52 lps2-601 ade2 soci62 6,810	A Adams
SERVER	MMT's lou2-3, IE2 anal-52 trp I-1 hini mys2-66	5 Brown

Cheung, K.-H., Samwald, M., Auerbach, R. K. & Gerstein, M. B. 2010. Structured digital tables on the Semantic Web: toward a structured digital literature. *Molecular Systems Biology, 6, 403*.

Health Informatics - Andreas Holzinger

69

OWL class constructors



Interse	ection/conjunctio	n of concepts,
	Speak: C1 and	Cn

Constructor	DL syntax	Example
Intersection	С1 ппСп	Anatomical_Abnormality Pathological_Function
Union	$C_1 \sqcup \ldots \sqcup C_n$	Body_Substance u Organic_Chemical
Complement	¬C	-Invertebrate
One of	$X_1 \sqcup \ldots \sqcup X_n$	Oestrogen u Progesterone
All values from	∀P.C	∀co_occurs_with.Plant
Some values	∃P.C	∃co_occurs_with.Animal
Max cardinality	$\leq nP$	1has_ingredient
Min cardinality	$\geq nP$	≥ 2. ingredient

Universal Restriction
Speak: All P-successors are in

Bhatt et al. (2009)

eak: An P-successor exists in C

Example for (2) Web Ontology Language OWL



DL = Description Logic		Concept inclusion, Speak: All C1 are C2		
Axiom Concept equivalence Speak: C1 is equivalent to C2	DL syntay	Example		
Sub class	$C_1 \sqsubseteq C_2$	Alga ⊑ Plant ⊑ Organism		
Equivalent class	$C_1 \equiv C_2$	Cancer = Neoplastic Process		
Disjoint with	$C_1 \sqsubseteq \neg C_2$	Vertebrate ¬Invertebrate		
Same individual	$x_1 \equiv x_2$	Blue_Shark = Prionace_Glauca		
Different from	$x_1 \sqsubseteq \neg x_2$	Sea Horse <u></u> ¬Horse		
Sub property	$P_1 \sqsubseteq P_2$	has_mother ⊑ has_parent		
Equivalent property	$P_1 \equiv P_2$	treated_by = cured_by		
Inverse	$P_1 \equiv P_2^-$	location_of = has_location		
Transitive property	$P^+ \sqsubseteq P$	part_of ⁺ part_of		
Functional property	$\top \sqsubseteq \leq 1P$	⊤ ⊑≤ 1has_tributary		
Inverse functional property	⊤ ⊑≤ 1 <i>P</i> −	⊤ ⊑ ≤ 1 has_scientific_name		

Bhatt, M., Rahayu, W., Soni, S. P. & Wouters, C. (2009) Ontology driven semantic profiling and retrieval in medical information systems. *Web Semantics: Science, Services and Agents on the World Wide Web, 7, 4, 317-331.*

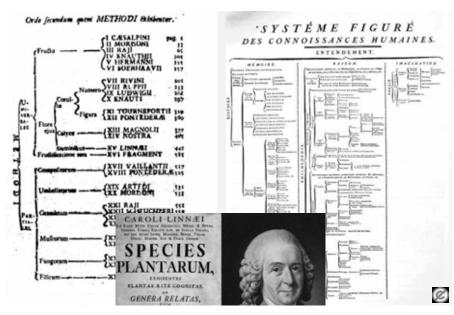
Health Informatics - Andreas Holzinger

70



06 Medical Classifications

Health Informatics – Andreas Holzinger 71 Health Informatics – Andreas Holzinger



 International Classification of Diseases (ICD) Systematized Nomenclature of Medicine (SNOMED)

approx. 100+ various classifications in use:

Since the classification by Carl von Linne (1735)

Medical Subject Headings (MeSH)

- Foundational Model of Anatomy (FMA)
- Gene Ontology (GO)
- Unified Medical Language System (UMLS)
- Logical Observation Identifiers Names & Codes (LOINC)
- National Cancer Institute Thesaurus (NCI Thesaurus)

Health Informatics - Andreas Holzinger

International Classification of Diseases (ICD)





http://www.who.int/classifications/icd/en

International Classification of Diseases (ICD)



- 1629 London Bills of Mortality
- 1855 William Farr (London, one founder of medical statistics): List of causes of death, list of diseases
- 1893 von Jacques Bertillot: List of causes of death
- 1900 International Statistical Institute (ISI) accepts Bertillot's list
- 1938 5th Edition
- 1948 WHO
- 1965 ICD-8
- 1989 ICD-10
- 2015 ICD-11 due
- 2018 ICD-11 adopt



Health Informatics - Andreas Holzinger



- 1965 SNOP, 1974 SNOMED, 1979 SNOMED II
- 1997 (Logical Observation Identifiers Names and Codes (LOINC) integrated into SNOMED
- 2000 SNOMED RT, 2002 SNOMED CT



SNOMED CT® Technical Reference Guide

January 2011 International Release (US English)

http://www.isb.nhs.uk/documents/isb-0034/amd-26-2006/techrefguid.pdf

Health Informatics - Andreas Holzinger

77

Medical Subject Headings (MeSH)



- MeSH thesaurus is produced by the National Library of Medicine (NLM) since 1960.
- Used for cataloging documents and related media and as an <u>index</u> to search these documents in a database and is part of the metathesaurus of the Unified Medical Language System (UMLS).
- This thesaurus originates from keyword lists of the Index Medicus (today Medline);
- MeSH thesaurus is polyhierarchic, i.e. every concept can occur multiple times. It consists of the three parts:
 - 1. MeSH Tree Structures,
 - 2. MeSH Annotated Alphabetic List and
 - 3. Permuted MeSH.

Α

24184005|Finding of increased blood pressure (finding) → 38936003|Abnormal blood pressure (finding) AND roleGroup SOME (363714003|Interprets (attribute) SOME 75367002|Blood pressure (observable entity))

Е

12763006|Finding of decreased blood pressure (finding) → 392570002|Blood pressure finding (finding) AND roleGroup SOME (363714003|Interprets (attribute) SOME 75367002|Blood pressure (observable entity))

Rector, A. L. & Brandt, S. (2008) Why Do It the Hard Way? The Case for an Expressive Description Logic for SNOMED. *Journal of the American Medical Informatics Association*, 15, 6, 744-751.

Health Informatics - Andreas Holzinger

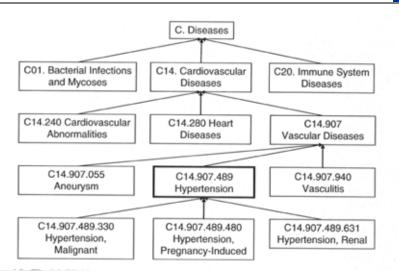
78

The 16 trees in MeSH



- 1. Anatomy [A]
- 2. Organisms [B]
- 3. Diseases [C]
- 4. Chemicals and Drugs [D]
- 5. Analytical, Diagnostic and Therapeutic Techniques and Equipment [E]
- 6. Psychiatry and Psychology [F]
- 7. Biological Sciences [G]
- 8. Natural Sciences [H]
- 9. Anthropology, Education, Sociology, Social Phenomena [I]
- 10. Technology, Industry, Agriculture [J]
- 11. Humanities [K]
- 12. Information Science [L]
- 13. Named Groups [M]
- 14. Health Care [N]
- 15. Publication Characteristics [V]
- 16. Geographicals [Z]

W/



Hersh, W. (2010) Information Retrieval: A Health and Biomedical Perspective. New York, Springer.

Health Informatics - Andreas Holzinger

81

MeSH Interactive Tree-Map Visualization (see L 9)





Eckert, K. (2008) A methodology for supervised automatic document annotation. *Bulletin of IEEE Technical Committee on Digital Libraries TCDL, 4, 2.*

National Library of Medicine - Medical Subject Headings

2011 MeSH

MeSH Descriptor Data

Between to Entry Priger

Standard View. Go to Concept View; Go to Expanded Concept View

MeSH Heading	Hypertension
Tree Number	C14.907.489
Annotation	not for intracrantal or intraocular pressure; relation to <u>BLOOD PRESSURE</u> : Manual 23.27; Goldblatt kidney is <u>HYPERTENSION</u> , GOLDBLATT see <u>HYPERTENSION</u> , <u>RENOVASCULAR</u> ; hypertension with kidney disease is probably <u>HYPERTENSION</u> , <u>RENAL</u> , not <u>HYPERTENSION</u> ; venous hypertension: Index under <u>VENOUS PRESSURE</u> (IM) & do not coordinate with <u>HYPERTENSION</u> ; <u>PREHYPERTENSION</u> is also available.
Scope Note	Persistently high systemic arterial BLOOD PRESSURE Based on multiple readings (_BLOOD PRESSURE DETERMINATION), hypertension is currently defined as when SYSTOLIC PRESSURE is consistently greater than 140 mm Hg or when DIASTOLIC PRESSURE is consistently 90 mm Hg or more.
Entry Term	Blood Pressure, High
See Also	Antihypertensive Agents
See Also	Vascular Resistance
Allowable Qualifiers	BL CF CLCL ON CO OH DLDT EC EH EM EN EP ET GE HLIM ME MLMO NU PA PC PP PS PX RA RH RERT SU TH UR US VE VI
Date of Entry	19990101
Unique ID	D006973

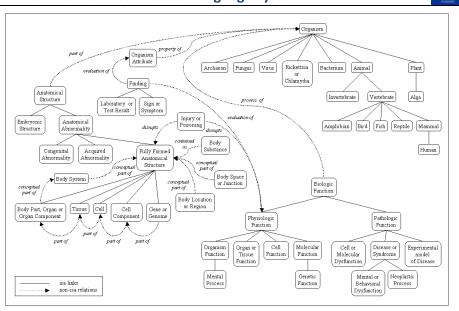
http://www.nlm.nih.gov/mesh/

Health Informatics - Andreas Holzinger

82

UMLS – Unified Medical Language System





http://www.nlm.nih.gov/research/umls/



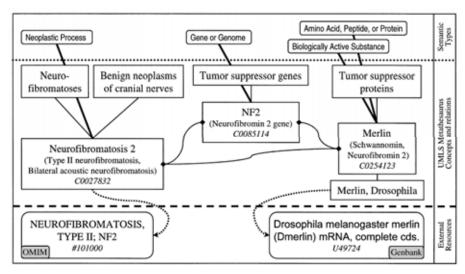


Health Informatics - Andreas Holzinger

85

Example of proteins and diseases in the UMLS

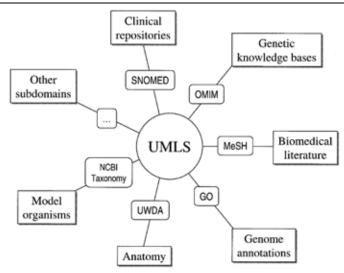




Bodenreider, O. (2004) The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Research*, *32*, *D267-D270*.

UMLS Metathesaurus integrates sub-domains





Bodenreider, O. (2004) The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Research*, 32, D267-D270.

Health Informatics - Andreas Holzinger

86



Conclusion and Future Challenges

W

- To find a trade-off between standardization and personalization [1];
- The large amounts of non-standardized data and unstructured information ("free text") [2];
- Low integration of standardized terminologies in the daily clinical practice (Who is using e.g. SNOMED, MeSH, UMLS in daily routine?);
- Low acceptance of classification codes amongst practitioners;
- Holmes, C., Mcdonald, F., Jones, M., Ozdemir, V., Graham, J. E. 2010. Standardization and Omics Science: Technical and Social Dimensions Are Inseparable and Demand Symmetrical Study. Omics-Journal of Integr. Biology, 14, (3), 327-332.
- Holzinger, A., Schantl, J., Schroettner, M., Seifert, C. & Verspoor, K. 2014. Biomedical Text Mining: State-of-the-Art, Open Problems and Future Challenges. In: LNCS 8401. Berlin Heidelberg: Springer pp. 271-300.

Health Informatics - Andreas Holzinger

89



Thank you!

- Data fusion Data integration in the life sciences
- Self learning stochastic ontologies [1]
- Interactive, integrative machine learning and interactive ontologies - human-in-the-loop
- Never ending learning machines [2] for automatically building knowledge spaces
- Integrating ontologies in daily work
- Knowledge and context awareness

[1] Ongenae, F., Claeys, M., Dupont, T., Kerckhove, W., Verhoeve, P., Dhaene, T. & De Turck, F. 2013. A probabilistic ontology-based platform for self-learning context-aware healthcare applications. Expert Systems with Applications, 40, (18), 7629-7646.

[2] Carlson, A., Betteridge, J., Kisiel, B., Settles, B., Hruschka Jr, E. R. & Mitchell, T. M. 2010. Toward an Architecture for Never-Ending Language Learning. Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence (AAAI-10). Atlanta: AAAI. 1306-1313.

Health Informatics - Andreas Holzinger

90



Questions



■ The Quiz-Slide will be shown during the course

