

Andreas Holzinger  
VO 709.049 Medical Informatics  
21.10.2015 11:15-12:45

## Lecture 02 Back to the Future – Fundamentals of Data, Information and Knowledge

a.holzinger@tugraz.at  
Tutor: markus.plass@student.tugraz.at  
Web: <http://hci-kdd.org/biomedical-informatics-big-data>

A. Holzinger 709.049 1/74 Med Informatics L02

Schedule

- 1. Introduction: Computer Science meets Life Sciences, challenges and future directions
- 2. Back to the future: Fundamentals of Data, Information and Knowledge
- 3. Structured Data: Coding, Classification (ICD, SNOMED, MeSH, UMLS)
- 4. Biomedical Databases: Acquisition, Storage, Information Retrieval and Use
- 5. Semi structured and weakly structured data (structural homologies)
- 6. Multimedia Data Mining and Knowledge Discovery
- 7. Knowledge and Decision: Cognitive Science & Human-Computer Interaction
- 8. Biomedical Decision Making: Reasoning and Decision Support
- 9. Intelligent Information Visualization and Visual Analytics
- 10. Biomedical Information Systems and Medical Knowledge Management
- 11. Biomedical Data: Privacy, Safety and Security
- 12. Methodology for Information Systems: System Design, Usability and Evaluation

A. Holzinger 709.049 2/74 Med Informatics L02

Keywords

- Computational space (high-dimensional)
- Data structures
- DIK-Model
- DIKW-Model
- Dimensionality of data
- Information complexity
- Information entropy
- Perceptual space (low-dimensional)
- Standardization versus Structurization

A. Holzinger 709.049 3/74 Med Informatics L02

Learning Goals

- ... be aware of the types and categories of different data sets in biomedical informatics;
- ... know some differences between data, information, knowledge and wisdom;
- ... be aware of standardized/non-standardized and well-structured/un-structured data;
- ... have a basic overview on information theory and the concept of information entropy;

A. Holzinger 709.049 4/74 Med Informatics L02

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 filtered according to their relevance 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);
- **Artifact/surrogate** = error or anomaly in the perception or representation of information through 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 organizing data to use it efficiently;
- **Deduction** = deriving a particular valid conclusion from a set of general premises;
- **DIK-Model** = Data-Information-Knowledge three level model
- **DIKW-Model** = Data-Information-Knowledge-Wisdom four 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;

A. Holzinger 709.049 5/74 Med Informatics L02

Advance Organizer (2/2)

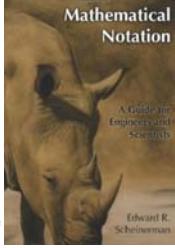
- **Induction** = deriving a likely general conclusion from a set of particular statements;
- **Information** = derived from the data by interpretation (with feedback to the clinician);
- **Information Entropy** = a measure for uncertainty: highly structured data contain low entropy, if everything is 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 multi-variate
- **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

A. Holzinger 709.049 6/74 Med Informatics L02

**Common Mathematical Notations with LaTeX commands**

*"In mathematics you don't understand things. You just get used to them" – John von Neumann*

<b>Data</b>	
$n$	Number of samples
$d$	Number of input variables
$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$	Matrix of input samples
$\mathbf{y} = [y_1, \dots, y_n]$	Vector of output samples
$\mathbf{Z} = [\mathbf{X}, \mathbf{y}]$	Combined input-output training data or
$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]$	Representation of data points in a feature space
<b>Distribution</b>	
$P$	Probability
$F(\mathbf{x})$	Cumulative probability distribution function (cdf)
$p(\mathbf{x})$	Probability density function (pdf)
$p(\mathbf{x}, \omega)$	Joint probability density function
$p(\omega \mathbf{x})$	Probability density function, which is parameterized
$p(\mathbf{y} \mathbf{x})$	Conditional density
$t(\mathbf{x})$	Target function



A. Holzinger 709.049      7/74      Med Informatics L02

**Glossary**

- ApEn = Approximate Entropy;
- $\mathbb{C}_{\text{data}}$  = Data in computational space;
- DIK = Data-Information-Knowledge-3-Level Model;
- DIKW = Data-Information-Knowledge-Wisdom-4-Level Model;
- GraphEn = Graph Entropy;
- H = Entropy (General);
- HRV = Heart Rate Variability;
- MaxEn = Maximum Entropy;
- MinEn = Minimum Entropy;
- NE = Normalized entropy (measures the relative informational content of both the signal and noise);
- $\mathbb{P}_{\text{data}}$  = Data in perceptual space;
- PDB = Protein Data Base;
- SampEn = Sample Entropy;

A. Holzinger 709.049      8/74      Med Informatics L02

**Key Problems**

- Heterogeneous, distributed, inconsistent data sources (need for **data integration & fusion**) [1]
- **Complex data** (high-dimensionality – challenge of dimensionality reduction and visualization) [2]
- Noisy, uncertain, missing, dirty, and imprecise, imbalanced data (challenge of **pre-processing**)
- The discrepancy between data-information-knowledge (**various definitions**)
- **Big data** sets (manual handling of the data is awkward, and often impossible) [3]

1. Holzinger A, Dehmer M, & Jurisica I (2014) Knowledge Discovery and Interactive Data Mining in Bioinformatics - State-of-the-Art, future challenges and research directions. *BMC Bioinformatics* 15(S6):1.  
2. Hund, M., Sturm, W., Schreck, T., Ullrich, T., Keim, D., Majaric, L., & Holzinger, A. 2015. Analysis of Patient Groups and Immunization Based on Subspace Clustering. In: *LNAI 9250*, 358-368.  
3. Holzinger, A., Stocker, C. & Dehmer, M. 2014. Big Complex Biomedical Data: Towards a Taxonomy of Data. In *CCIS 455*. Springer. 3-18.

A. Holzinger 709.049      9/74      Med Informatics L02

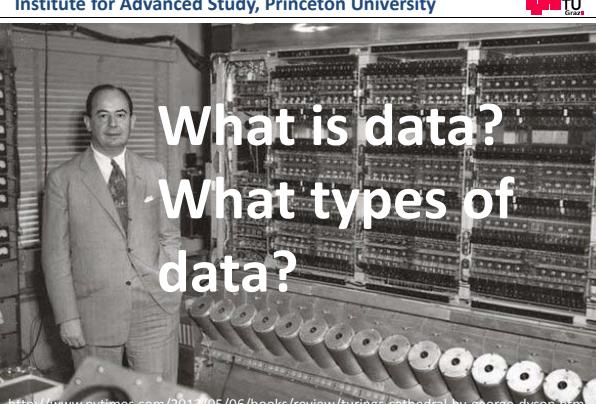
**Traditional Statistics versus Machine Learning**

<ul style="list-style-type: none"> <li>▪ Data in traditional Statistics</li> <li>▪ Low-dimensional data (<math>&lt; \mathbb{R}^{100}</math>)</li> <li>▪ Problem: Much noise in the data</li> <li>▪ Not much structure in the data but it can be represented by a simple model</li> </ul>	<ul style="list-style-type: none"> <li>▪ Data in Machine Learning</li> <li>▪ High-dimensional data (<math>\gg \mathbb{R}^{100}</math>)</li> <li>▪ Problem: not noise, but complexity</li> <li>▪ Much structure, but the structure but can <b>not</b> be represented by a simple model</li> </ul>
--	--

Lecun, Y., Bengio, Y. & Hinton, G. 2015. Deep learning. *Nature*, 521, (7553), 436-444.

A. Holzinger 709.049      10/74      Med Informatics L02

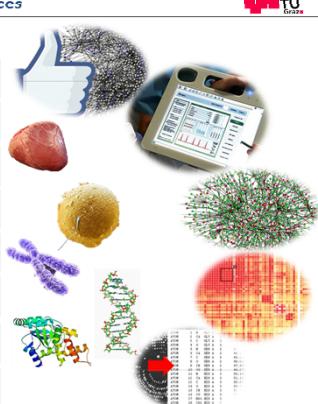
**Institute for Advanced Study, Princeton University**



<http://www.nytimes.com/2012/05/06/books/review/turings-cathedral-by-george-dyson.html>

A. Holzinger 709.049      11/74      Med Informatics L02

**Slide 2-1: Biomedical Data Sources**

<p><math>10^{-12}</math></p> <ul style="list-style-type: none"> <li><b>Collective</b> </li> <li><b>Individual</b> </li> <li><b>Tissue</b> </li> <li><b>Cell</b> </li> <li><b>Bacteria</b> </li> <li><b>Virus</b> </li> <li><b>Molecule</b> </li> <li><b>Atom</b> </li> </ul>	
--	---

A. Holzinger 709.049      12/74      Med Informatics L02

**Slide 2-2: Taxonomy of data**



- **Physical level** -> bit = binary digit = **basic** indissoluble unit (= Shannon, Sh), ≠ Bit (!) in Quantum Systems -> qubit
- **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, ...)

**Slide 2-3: Example Data Structures (1/3): List**

The diagram illustrates two examples of list data structures:

- TYPE link = obj{ node : 1, next : 2, key : 3, linkType : 4 }**: A linked list node structure.
- VAR p, q : link :=**: Declaration of pointers p and q.
- p := NEW(link);**: Initialization of pointer p to a new list node.
- q := NEW(link);**: Initialization of pointer q to a new list node.
- p^.key := 1;**: Setting the key value of node p to 1.
- q^.key := 2;**: Setting the key value of node q to 2.
- q := p^.link();**: Setting pointer q to the next node of p.
- p := q^.link();**: Setting pointer p to the next node of q.
- q := p^.link();**: Setting pointer q to the next node of p.

The diagram also shows a CAP-DNA Complex structure with labels for the CAP protein, DNA, Turn-Helix, and B-DNA.

**B CAPrecognition via DNALogo**

A sequence logo for the sequence TGTGAACACAT is shown, with a background frequency of 0.2. The logo is composed of four columns representing A, T, C, and G, with heights indicating their relative frequencies at each position.

**C CAP Heits-Turn-Helix Logo**

A sequence logo for the sequence RGAAGCTTTCAT is shown, with a background frequency of 0.2. The logo is composed of four columns representing R, F, Y, and S, with heights indicating their relative frequencies at each position.

## Slide 2-4: 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.

$$W = \begin{bmatrix} 0 & w_{12} & w_{13} & 0 & 0 \\ 0 & 0 & w_{23} & w_{24} & 0 \\ w_{31} & 0 & 0 & 0 & w_{35} \\ 0 & w_{42} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{54} & 0 \end{bmatrix}$$

A. Holzinger 709.049

15/74

Med Informatics L02

# Data Integration and Data Fusion in the Life Sciences

## Slide 2-6: “Big Data” pools in the health domain

**Slide 2-7a: Omics-data integration (1/2)**

Genomics	Transcriptomics	Proteomics	Metabolomics	Protein-DNA interactions	Protein-protein interactions	Fluxomics	Phenomics
Genomics (sequence annotation)	Transcriptomics (microarray, SAGE)	Proteomics (abundance, post-translational modification)	Metabolomics (metabolic abundance)	Protein-DNA interactions (ChIP-chip)	Protein-protein interactions (yeast 2H, co-IP-MS)	Fluxomics (isotopic tracing)	Phenomics (phenotype arrays, RNA screens, synthetic lethals)
• ORF validation • Regulatory element identification <sup>[1]</sup>	• SNP effect on protein activity or abundance • Protein transcript correlation <sup>[2]</sup>	• Enzyme annotation	• Gene regulatory networks <sup>[3]</sup>	• Binding-site identification <sup>[4]</sup>	• Functional annotation <sup>[5]</sup>	• Functional annotation <sup>[6]</sup>	• Functional annotation <sup>[7]</sup>
Transcriptomics (microarray, SAGE)	Proteomics (abundance, post-translational modification)	Metabolomics (metabolic abundance)	Protein-DNA interactions (ChIP-chip)	Protein-protein interactions (yeast 2H, co-IP-MS)	Fluxomics (isotopic tracing)	Phenomics (phenotype arrays, RNA screens, synthetic lethals)	Phenomics (phenotype arrays, RNA screens, synthetic lethals)
• ORF validation • Regulatory element identification <sup>[1]</sup>	• SNP effect on protein activity or abundance • Protein transcript correlation <sup>[2]</sup>	• Enzyme annotation	• Gene regulatory networks <sup>[3]</sup>	• Binding-site identification <sup>[4]</sup>	• Functional annotation <sup>[5]</sup>	• Functional annotation <sup>[6]</sup>	• Functional annotation <sup>[7]</sup>
Genomics (sequence annotation)	Transcriptomics (microarray, SAGE)	Proteomics (abundance, post-translational modification)	Metabolomics (metabolic abundance)	Protein-DNA interactions (ChIP-chip)	Protein-protein interactions (yeast 2H, co-IP-MS)	Fluxomics (isotopic tracing)	Phenomics (phenotype arrays, RNA screens, synthetic lethals)

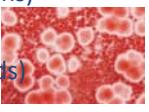
Joyce, A. R. & Palsson, B. Ø. 2006. The model organism as a system: integrating 'omics' data sets. *Nature Reviews Molecular Cell Biology*, 7, 198-210.

A. Holzinger 709.049 19/74 Med Informatics L02

**Slide 2-7b: -Omics-data integration (2/2)**

- Genomics (sequence annotation)
- Transcriptomics (microarray)
- Proteomics (Proteome Databases)
- Metabolomics (enzyme annotation)
- Fluxomics (isotopic tracing, metabolic pathways)
- Phenomics (biomarkers)
- Epigenomics (epigenetic modifications)
- Microbiomics (microorganisms)
- Lipidomics (pathways of cellular lipids)

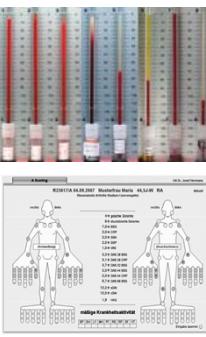




A. Holzinger 709.049 20/74 Med Informatics L02

**Slide 2-8: Example of typical clinical data sets**

- 50+ Patients per day ~ 5000 data points per day ...
- Aggregated with specific scores (Disease Activity Score, DAS)
- Current patient status is related to previous data
- = convolution over time
- ⇒ time-series data



Simonic, K. M., Holzinger, A., Bloice, M. & Hermann, J. (2011). Optimizing Long-Term Treatment of Rheumatoid Arthritis with Systematic Documentation. *Pervasive Health - 5th International Conference on Pervasive Computing Technologies for Healthcare, Dublin, IEEE*, 550-554.

A. Holzinger 709.049 21/74 Med Informatics L02

**Slide 2-9: Standardization vs. Structurization**

Weakly-Structured	Omics Data	Natural Language Text
Well-Structured	Databases	Libraries
	RDF, OWL	
Standardized		Non-Standardized

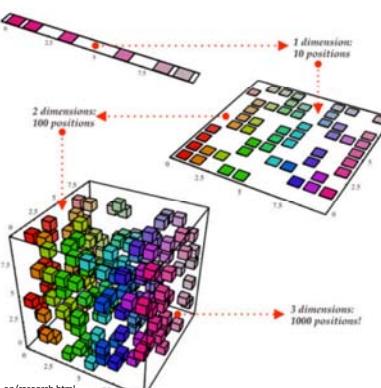
Holzinger, A. (2011) Weakly Structured Data in Health-Informatics: The Challenge for Human-Computer Interaction. In: Bagheri, N., Baxter, G., Dow, L. & Kimani, S. (Eds.) *Proceedings of INTERACT 2011 Workshop: Promoting and supporting healthy living by design*. Lisbon, IFIP, 5-7.

A. Holzinger 709.049 22/74 Med Informatics L02

**Note: The curse of dimensionality**

Bengio, S. & Bengio, Y. 2000. Taking on the curse of dimensionality in joint distributions using neural networks. *IEEE Transactions on Neural Networks*, 11, (3), 550-557.

[http://www.iro.umontreal.ca/~bengioy/yoshua\\_en/research.html](http://www.iro.umontreal.ca/~bengioy/yoshua_en/research.html)



A. Holzinger 709.049 23/74 Med Informatics L02

**Slide 2-10: Data Dimensionality examples**

- 0-D data = a data point existing isolated from other data, e.g. integers, letters, Booleans, etc.
- 1-D data = consist of a string of 0-D data, e.g. Sequences representing nucleotide bases and amino acids, SMILES etc.
- 2-D data = having spatial component, such as images, NMR-spectra etc.
- 2.5-D data = can be stored as a 2-D matrix, but can represent biological entities in three or more dimensions, e.g. PDB records
- 3-D data = having 3-D spatial component, e.g. image voxels, e-density maps, etc.
- H-D Data = data having arbitrarily high dimensions

A. Holzinger 709.049 24/74 Med Informatics L02

**Example: 1-D data (univariate sequential data objects)**

SMILES (Simplified Molecular Input Line Entry Specification)

... is a compact machine and human-readable chemical nomenclature:

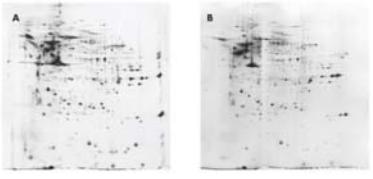
e.g. Viagra:  
CCc1nn(C)c2c(=O)[nH]c(nc12)c3cc(ccc3OCC)S(=O)(=O)N4CCN(C)CC4

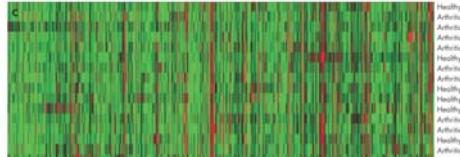
...is Canonicalizable  
 ...is Comprehensive  
 ...is Well Documented

[http://www.daylight.com/dayhtml\\_tutorials/languages/smiles/index.html](http://www.daylight.com/dayhtml_tutorials/languages/smiles/index.html)

A. Holzinger 709.049 25/74 Med Informatics L02

**Example: 2-D data (bivariate data)**

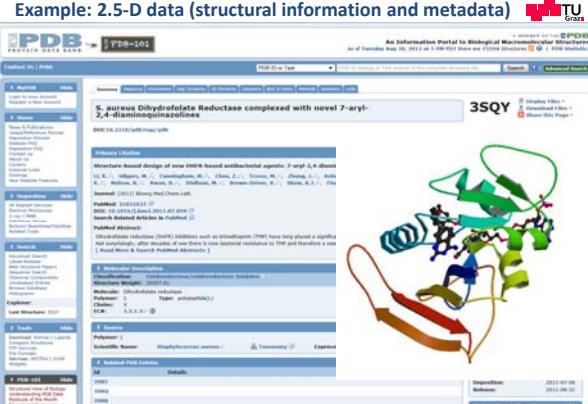




Kastrinaki et al. (2008) Functional, molecular & proteomic characterisation of bone marrow mesenchymal stem cells in rheumatoid arthritis. *Annals of Rheumatic Diseases*, 67, 6, 741-749.

A. Holzinger 709.049 26/74 Med Informatics L02

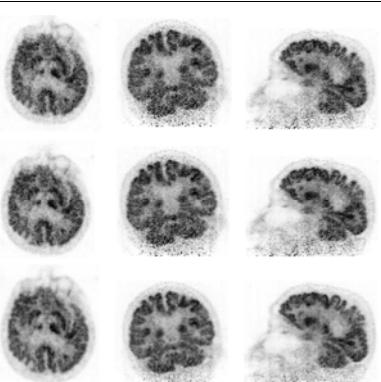
**Example: 2.5-D data (structural information and metadata)**



<http://www.pdb.org>

A. Holzinger 709.049 27/74 Med Informatics L02

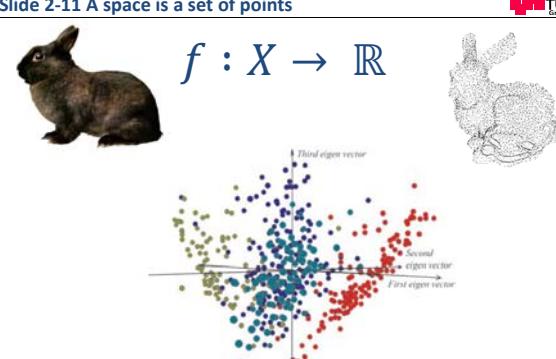
**Example: 3-D Voxel data (volumetric picture elements)**



Scheins, J. J., Herzog, H. & Shah, N. J. (2011) Fully-3D PET Image Reconstruction Using Scanner-Independent, Adaptive Projection Data and Highly Rotation-Symmetric Voxel Assemblies. *Medical Imaging, IEEE Transactions on*, 30, 3, 879-892.

A. Holzinger 709.049 28/74 Med Informatics L02

**Slide 2-11 A space is a set of points**



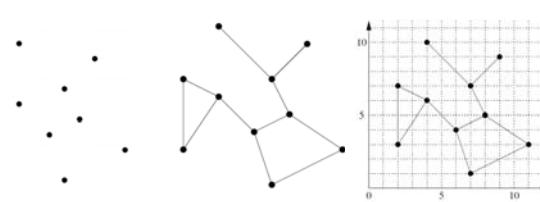
$f : X \rightarrow \mathbb{R}$

Hou, J., Sims, G. E., Zhang, C. & Kim, S.-H. 2003. A global representation of the protein fold space. *Proceedings of the National Academy of Sciences*, 100, (5), 2386-2390.

A. Holzinger 709.049 29/74 Med Informatics L02

**Slide 2-12 Point Cloud Data Sets**

Let us collect  $n$ -dimensional  $i$  observations:  $x_i = [x_{i1}, \dots, x_{in}]$



Point cloud in  $\mathbb{R}^2$       topological space      metric space

Zomorodian, A. J. 2005. *Topology for computing*, Cambridge (MA), Cambridge University Press.

A. Holzinger 709.049 30/74 Med Informatics L02

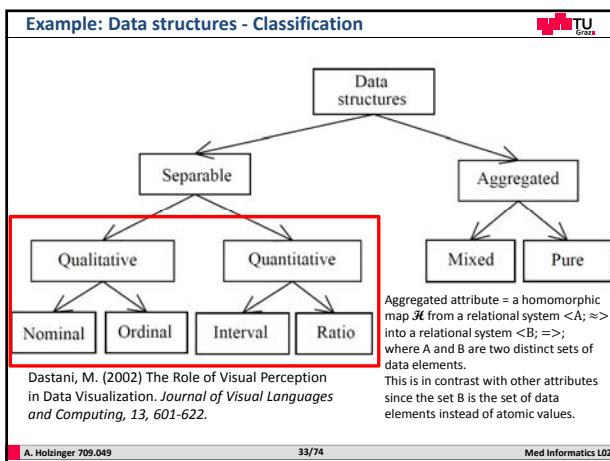
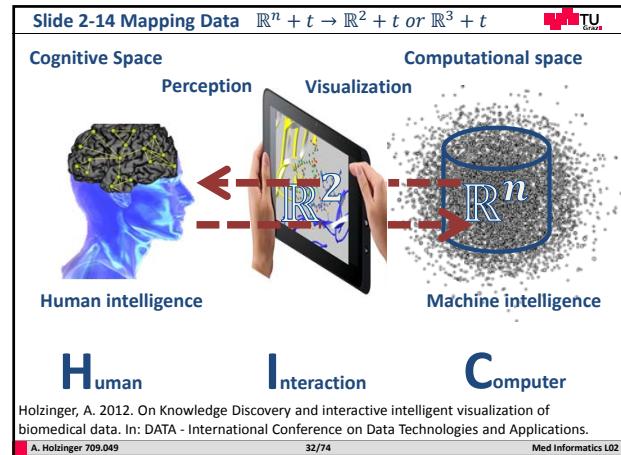
**Slide 2-13: Example Metric Space**

A set  $S$  with a metric function  $d$  is a metric space

$$d_{ij} = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2}$$

Doob, J. L. 1994. *Measure theory*, Springer New York.

A. Holzinger 709.049 31/74 Med Informatics L02

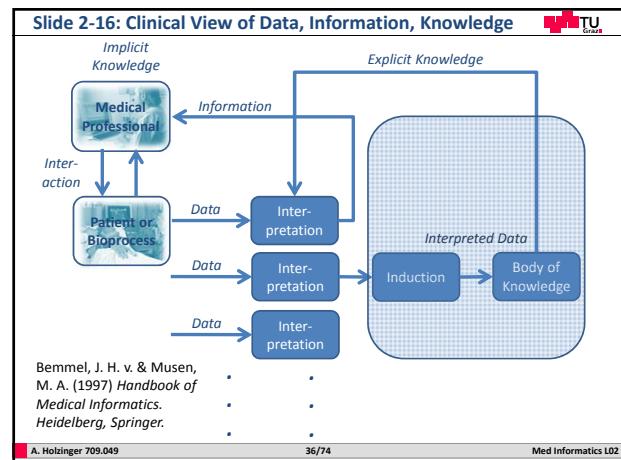


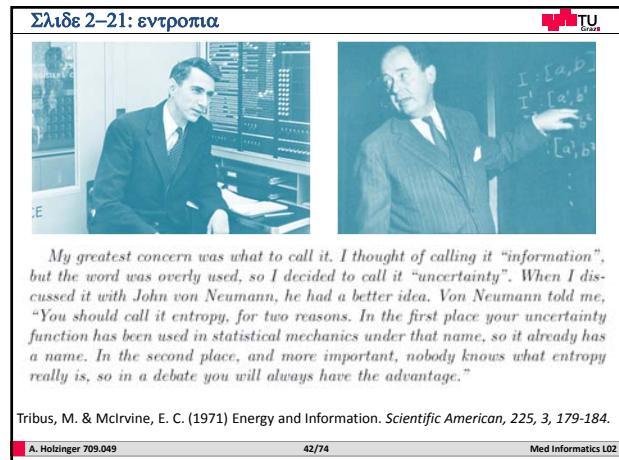
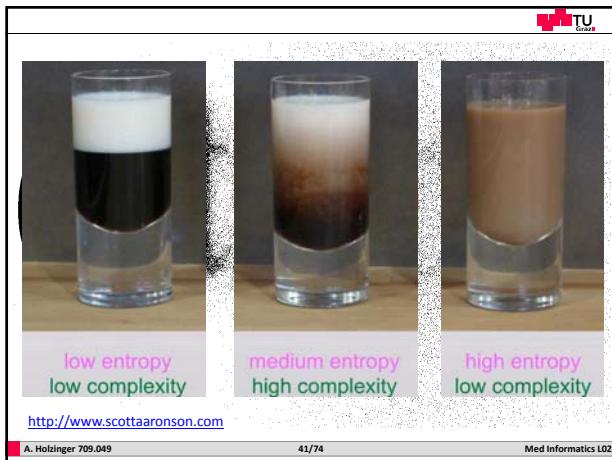
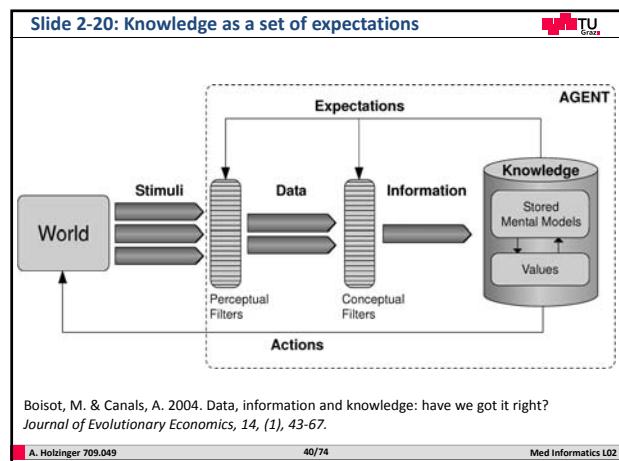
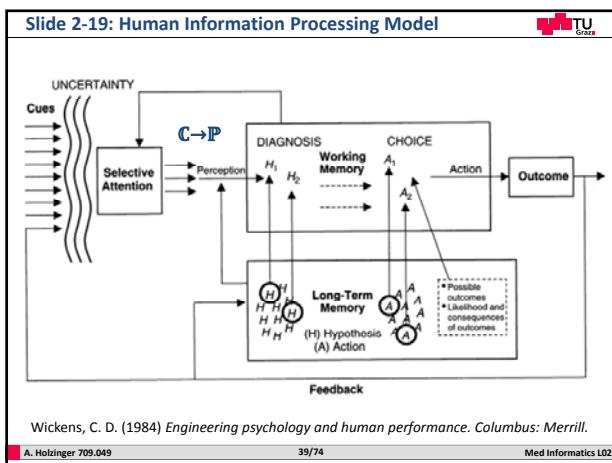
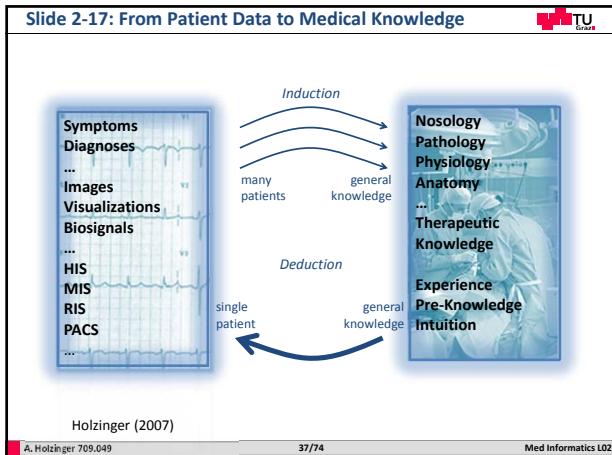
**Slide 2-15: Categorization of Data (Classic "scales")**

Scale	Empirical Operation	Mathem. Group Structure	Transf. in $\mathbb{R}$	Basic Statistics	Mathematical Operations
NOMINAL	Determination of equality	Permutation $x' = f(x)$ $x \dots 1\text{-to-}1$	$x \mapsto f(x)$	Mode, contingency correlation	$=, \neq$
ORDINAL	Determination of more/less	Isotonic $x' = f(x)$ $x \dots \text{mono-}\text{tonic incr.}$	$x \mapsto f(x)$	Median, Percentiles	$=, \neq, >, <$
INTERVAL	Determination of equality of intervals or differences	General linear $x' = ax + b$	$x \mapsto rx+s$	Mean, Std.Dev. Rank-Order Corr., Prod.-Moment Corr.	$=, \neq, >, <, +, -, \times$
RATIO	Determination of equality or ratios	Similarity $x' = ax$	$x \mapsto rx$	Coefficient of variation	$=, \neq, >, <, +, -, \times, \div$

Stevens, S. S. (1946) On the theory of scales of measurement. *Science*, 103, 677-680.

A. Holzinger 709.049 34/74 Med Informatics L02





**Slide 2-22: Entropy H as a measure for uncertainty (1/3)**

$Q \dots P = \{p_1, \dots, p_n\}$        $H(Q) = - \sum_{i=1}^n (p_i * \log p_i)$

$Qb = \{a_1, a_2\}$  with  $P = \{p, 1-p\}$

$H(Qb) = p * \log \frac{1}{p} + p * \log \frac{1}{1-p}$

Shannon, C. E. (1948) A Mathematical Theory of Communication. *Bell System Technical Journal*, 27, 379-423.

Shannon, C. E. & Weaver, W. (1949) *The Mathematical Theory of Communication*. Urbana (IL), University of Illinois Press.

A. Holzinger 709.049      43/74      Med Informatics L02

**Slide 2-23: A measure for uncertainty (2/3)**

$\log_2 \frac{1}{p} = -\log_2 p$

$H = - \sum_{i=1}^N p_i \log_2(p_i)$

Shannon, C. E. (1948) A Mathematical Theory of Communication. *Bell System Technical Journal*, 27, 379-423.

A. Holzinger 709.049      44/74      Med Informatics L02

**Slide 2-24: Entropy H as a measure for uncertainty (3/3)**

$H_B = - \sum_{k=1}^B p_k \log_2 p_k = -1 * \log_2(1) = 0$

$H_B = - \sum_{k=1}^B \frac{1}{B} \log_2 \frac{1}{B} = \log_2(B)$

$H = H_{min} = 0 \quad H = H_{max} = \log_2 N$

A. Holzinger 709.049      45/74      Med Informatics L02

**Entropic methods – what for?**

- 1) Set of noisy, complex data
- 2) Extract information out of the data
- 3) to support a previous set hypothesis
- Information + Statistics + Inference
- = powerful methods for many sciences
- Application e.g. in biomedical informatics for analysis of ECG, MRI, CT, PET, sequences and proteins, DNA, topography, and for modeling etc.;

Mayer, C., Bachler, M., Hortenhuber, M., Stocker, C., Holzinger, A. & Wassertheurer, S. 2014. Selection of entropy-measure parameters for knowledge discovery in heart rate variability data. *BMC Bioinformatics*, 15, (Suppl 6), S2.

A. Holzinger 709.049      46/74      Med Informatics L02

**Slide 2-25: An overview on the History of Entropy**

```

graph TD
    Bernoulli["Bernoulli (1713)  
Principle of Insufficient Reason"] --> Bayes["Bayes (1763), Laplace (1770)  
How to calculate the state of a system with a limited number of expectation values"]
    Maxwell["Maxwell (1859), Boltzmann (1871),  
Gibbs (1902) Statistical Modeling of problems in physics"] --> Shannon["Shannon (1948)  
Information Theory"]
    Pearson["Pearson (1900)  
Goodness of Fit measure"] --> Fisher["Fisher (1922)  
Maximum Likelihood"]
    Bayes --> JeffreysCox["Jeffreys, Cox (1939-1948)  
Statistical Inference"]
    Maxwell --> Shannon
    Pearson --> Fisher
    Fisher --> Shannon
    Shannon --> BayesianStatistics["Bayesian Statistics"]
    Shannon --> EntropyMethods["Entropy Methods"]
    Shannon --> GeneralizedEntropy["Generalized Entropy"]
    BayesianStatistics --> SeeNextSlide["See next slide"]
    EntropyMethods --> SeeNextSlide
    GeneralizedEntropy --> SeeNextSlide
  
```

confer also with: Golan, A. (2008) Information and Entropy Econometric: A Review and Synthesis. *Foundations and Trends in Econometrics*, 2, 1-2, 1-145.

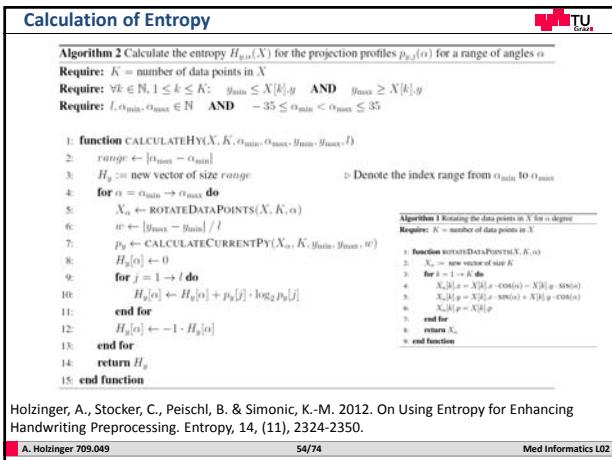
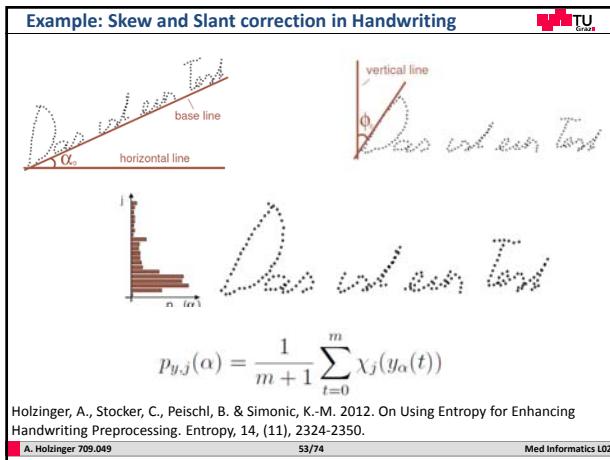
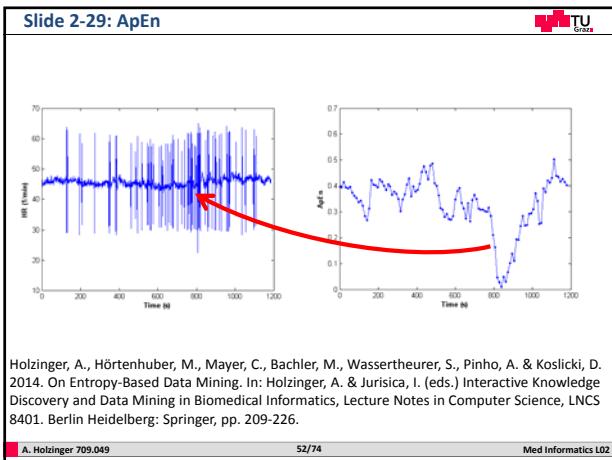
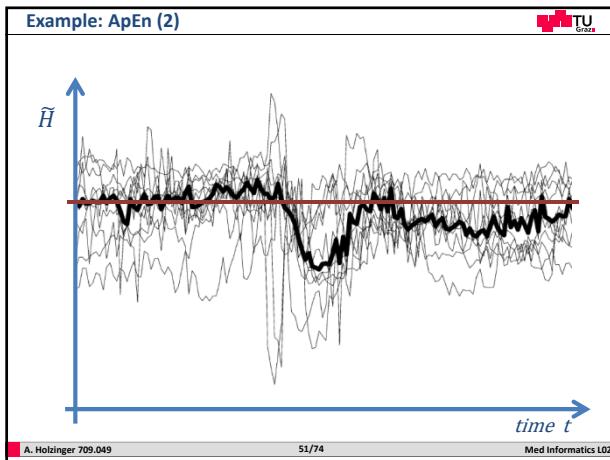
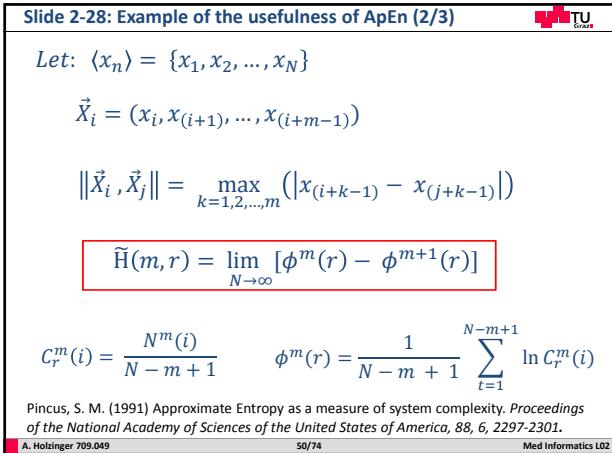
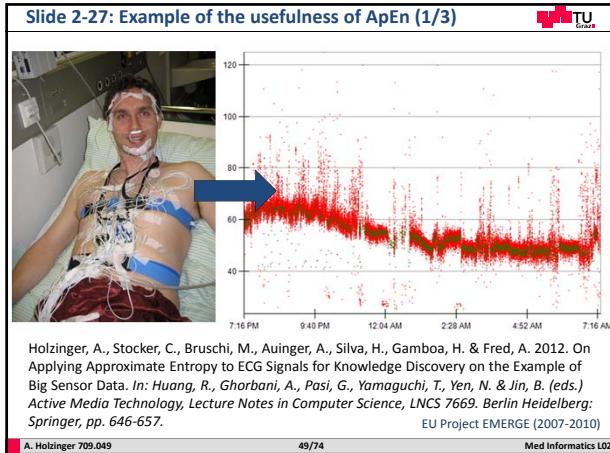
A. Holzinger 709.049      47/74      Med Informatics L02

**Slide 2-26: Towards a Taxonomy of Entropic Methods**

<p>Entropic Methods</p> <p>Jaynes (1957) <b>Maximum Entropy (MaxEn)</b></p> <p>Adler et al. (1965) <b>Topology Entropy (TopEn)</b></p> <p>Pincus (1991) <b>Approximate Entropy (ApEn)</b></p> <p>Richman (2000) <b>Sample Entropy (SamEn)</b></p>	<p>Generalized Entropy</p> <p>Renyi (1961) <b>Renyi-Entropy</b></p> <p>Mowshowitz (1968) <b>Graph Entropy (MinEn)</b></p> <p>Posner (1975) <b>Minimum Entropy (MinEn)</b></p> <p>Rubinstein (1997) <b>Cross Entropy (CE)</b></p> <p>Tsallis (1980) <b>Tsallis-Entropy</b></p>
---	---

Holzinger, A., Hörenhuber, M., Mayer, C., Bachler, M., Wassertheurer, S., Pinho, A. & Koslicki, D. 2014. On Entropy-Based Data Mining. In: Holzinger, A. & Jurisica, I. (eds.) *Lecture Notes in Computer Science*, LNCS 8401. Berlin Heidelberg: Springer, pp. 209-226.

A. Holzinger 709.049      48/74      Med Informatics L02



**Conclusion**

**H ...**

- ... is **robust** against noise;
- ... can be applied to **complex time series** with good replication;
- ... is **finite** for stochastic, noisy, composite processes;
- ... the values correspond directly to irregularities – good for detecting **anomalies**

A. Holzinger 709.049 55/74 Med Informatics L02



# Thank you!

A. Holzinger 709.049 56/74 Med Informatics L02

**Sample Questions (1)**

- Why is modeling of artifacts a huge problem?
- What do we need to transfer information into Knowledge?
- What type of data does the PDB basically store?
- What is the “curse of dimensionality”?
- What type of separable data is blood sedimentation rate?
- Is the mathematical operation “multiplication” allowed with ordinal data?
- What characterizes standardized data?
- Why are structural homologies interesting?
- How did Bemmel & van Mussen describe the clinical view on data, information and knowledge?
- Where are the differences between patient data and medical knowledge from a clinical viewpoint?
- Which weaknesses of the DIKW Model do you recognize?
- How do we get theories?
- What is the main limitation of transferring data from the computational space into the perceptual space from the viewpoint of the human information processing model?

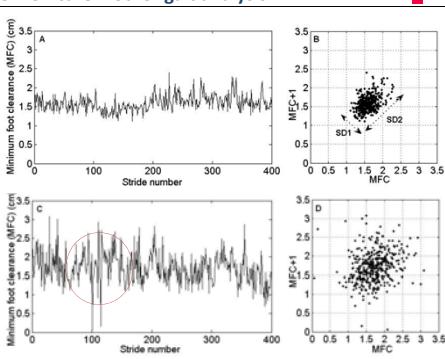
A. Holzinger 709.049 57/74 Med Informatics L02

**Sample Questions (2)**

- Why is the knowledge about human information processing necessary for medical informatics?
- What is the difference between the perceptual space and the computational space in terms of data, information and knowledge?
- What does information interaction mean?
- How does knowledge-assisted visualization work in principle?
- Why is non-structured data an rather incorrect term?
- Give an example of the data structure tree in biomedical informatics!
- Why is data quality important? What are the related issues?
- How do you ensure data accessibility?
- What is the main idea of Shannon’s Entropy?
- Why is Entropy interesting for medical informatics?
- What are typical entropic methods?
- What is the main purpose of Approximate Entropy?
- What is the big advantage of entropic methods?
- What are the differences of ApEn and SampEn?
- Which possibilities do you have with Graph Entropy Measures?

A. Holzinger 709.049 58/74 Med Informatics L02

**Back-up Slide: Poincare Plot for gait analysis**



Khandoker, A., Palaniswami, M. & Begg, R. (2008) A comparative study on approximate entropy measure and poincare plot indexes of minimum foot clearance variability in the elderly during walking. *Journal of NeuroEngineering and Rehabilitation*, 5, 1, 4.

A. Holzinger 709.049 59/74 Med Informatics L02

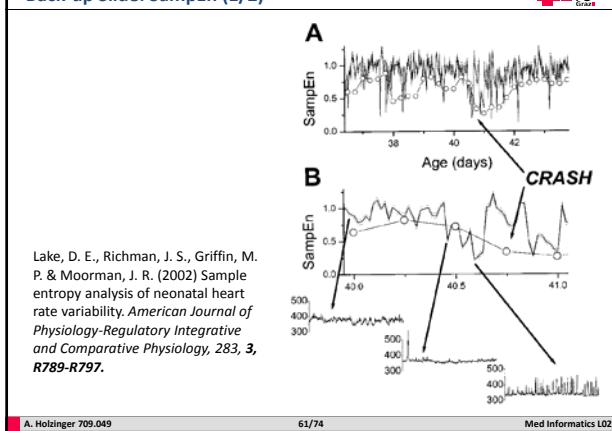
**Sample Exam Questions – Yes/No Answers**

01	An array is a composite data type on physical level.	<input type="checkbox"/> Yes	2 total
02	In a Von-Neumann machine “List” is a widely used data structure for applications which do not need random access.	<input type="checkbox"/> No	2 total
03	The edges in a graph can be multidimensional objects, e.g. vectors containing the results of multiple Gen-expression measures.	<input type="checkbox"/> Yes	2 total
04	Each item of data is composed of variables, and if such a data item is defined by more than one variable it is called a multivariable data item	<input type="checkbox"/> Yes	2 total
05	A dendrogram is a tree diagram frequently used to illustrate the arrangement of the clusters produced by hierarchical clustering.	<input type="checkbox"/> No	2 total
06	Nominal and ordinal data are parametric, and do assume a particular distribution.	<input type="checkbox"/> Yes	2 total
07	Abstraction is characterized by a cyclical process of generating possible explanations and testing those explanations.	<input type="checkbox"/> Yes	2 total
08	A metric space has an associated metric, which enables us to measure distances between points in that space and, in turn, implicitly define their neighborhoods.	<input type="checkbox"/> Yes	2 total
09	Induction consists of deriving a likely general conclusion from a set of particular statements.	<input type="checkbox"/> No	2 total
10	In the model of Boisot & Canale (2004), the perceptual filter orients the senses (e.g. visual sense) to certain types of stimuli within a certain physical range.	<input type="checkbox"/> Yes	2 total

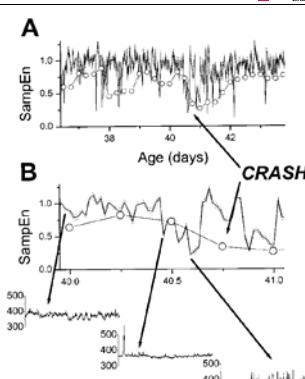
Sum of Question Block A (max. 20 points)

A. Holzinger 709.049 60/74 Med Informatics L02

## Backup Slide: SampEn (1/2)



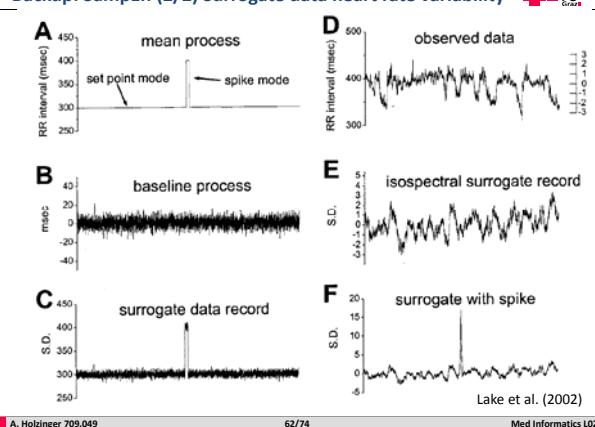
A. Holzinger 709.049



61/74

Med Informatics L02

## Backup: SampEn (2/2) Surrogate data heart rate variability



A. Holzinger 709.049

62/74

Med Informatics L02

## Backup Slide: Comparison ApEn - SampEn

ApEn

Given a signal  $x(n)=x(1), x(2), \dots, x(N)$ , where  $N$  is the total number of data points, ApEn algorithm can be summarized as follows [5]:

- 1) Form  $m$ -vectors,  $X(t)$  to  $X(N-m+1)$  defined by:  
 $X(t) = [x(t), x(t+1), \dots, X(t+m-1)] \quad t=1, N-m+1$  (1)
- 2) Define the distance  $d(X(t), X(j))$  between vectors  $X(t)$  and  $X(j)$  as the maximum absolute difference between their respective scalar components:  
 $d_m[X(t), X(j)] = \max_{k=0, m-1} \|x(t+k) - x(j+k)\|$  (2)
- 3) Define for each  $i$ , for  $i=1, N-m+1$ , let  
 $C_r^m(i) = l^m(i)/(N-m+1)$   
where  $l^m(i) = \text{no. of } d[X(t), X(i)] \leq r$  (3)
- 4) Take the natural logarithm of each  $C_r^m(i)$ , and average it over  $i$  as defined in step 3:  
 $\phi^m(r) = \frac{1}{N-m+1} \sum_{i=1}^{N-m+1} \ln(C_r^m(i))$  (4)
- 5) Increase the dimension to  $m+1$  and repeat steps 1) to 4).
- 6) Calculate ApEn value for a finite data length of  $N$ :  
 $ApEn(m, r, N) = \phi^m(r) - \phi^{m+1}(r)$  (5)

Xinnian, C. et al. (2005). Comparison of the Use of Approximate Entropy and Sample Entropy: Applications to Respiratory Signal. *Engineering in Medicine and Biology IEEE-EMBS 2005*, 4212-4215.

A. Holzinger 709.049

SampEn  
Given a signal  $x(n)=x(1), x(2), \dots, x(N)$ , where  $N$  is the total number of data points, SampEn algorithm can be summarized as follows [5]:

- 1) Form  $m$ -vectors,  $X(t)$  to  $X(N-m+1)$  defined by:  
 $X(t) = [x(t), x(t+1), \dots, X(t+m-1)] \quad t=1, N-m+1$  (6)
- 2) Define the distance  $d_m[X(t), X(j)]$  between vectors  $X(t)$  and  $X(j)$  as the maximum absolute difference between their respective scalar components:  
 $d_m[X(t), X(j)] = \max_{k=0, m-1} \|x(t+k) - x(j+k)\|$  (7)
- 3) Define for each  $i$ , for  $i=1, N-m$ , let  
 $B_r^m(r) = \frac{1}{N-m-1} \times \text{no. of } d_m[X(t), X(j)] \leq r, i \neq j$  (8)
- 4) Similarly, define for each  $i$ , for  $i=1, N-m$ , let  
 $A_r^m(r) = \frac{1}{N-m-1} \times \text{no. of } d_{m+1}[X(t), X(j)] \leq r, i \neq j$  (9)
- 5) Define  $B''(r) = \frac{1}{N-m} \sum_{i=1}^{N-m} B_r^m(r)$  (10)
- 6) Define  $A''(r) = \frac{1}{N-m} \sum_{i=1}^{N-m} A_r^m(r)$  (11)
- 7) SampEn value for a finite data length of  $N$  can be estimated:  
 $SampEn(m, r, N) = -\ln(A''(r)/B''(r))$  (12)

Dehmer, M. & Mowshowitz, A. (2011) A history of graph entropy measures. *Information Sciences*, 181, 1, 57-78.

A. Holzinger 709.049

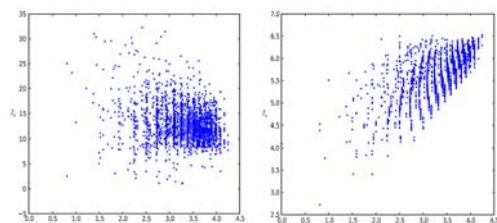
63/74

Med Informatics L02

## Backup Slide: Graph Entropy Measures

■ The most important question: Which kind of structural information does the entropy measure detect?

■ the topological complexity of a molecular graph is characterized by its number of vertices and edges, branching, cyclicity etc.



Dehmer, M. & Mowshowitz, A. (2011) A history of graph entropy measures. *Information Sciences*, 181, 1, 57-78.

A. Holzinger 709.049

64/74

Med Informatics L02

## Backup: English/German Subject Codes OEFOS 2012

106005	Bioinformatics	Bioinformatik
106007	Biostatistics	Biostatistik
304005	Medical Biotechnology	Medizinische Biotechnologie
305901	Computer-aided diagnosis and therapy	Computerunterstützte Diagnose und Therapie
304003	Genetic engineering, - technology	Gentechnik, -technologie
3906 (old)	Medical computer sciences	Medizinische Computerwissenschaften
305906	Medical cybernetics	Medizinische Kybernetik
305904	Medical documentation	Medizinische Dokumentation
305905	Medical informatics	Medizinische Informatik
305907	Medical statistics	Medizinische Statistik

<http://www.statistik.at>

A. Holzinger 709.049

65/74

Med Informatics L02

## Backup: English/German Subject Codes OEFOS 2012

102001	Artificial Intelligence	Künstliche Intelligenz
102032	Computational Intelligence	Computational Intelligence
102033	Data Mining	Data Mining
102013	Human-Computer Interaction	Human-Computer Interaction
102014	Information design	Informationsdesign
102015	Information systems	Informationssysteme
102028	Knowledge engineering	Knowledge Engineering
102019	Machine Learning	Maschinelles Lernen
102020	Medical Informatics	Medizinische Informatik
102021	Pervasive Computing	Pervasive Computing
102022	Software development	Softwareentwicklung
102027	Web engineering	Web Engineering

<http://www.statistik.at>

A. Holzinger 709.049

66/74

Med Informatics L02

**Backup Slide: Statistical Analysis Software (SAS)**

The slide shows a screenshot of the SAS website. The main navigation bar includes 'Support', 'Knowledge Base', 'Report', 'Training & Books', and 'Community'. Below this, there's a 'PDF' link and a 'Contents' section listing various SAS procedures. A specific section for the 'ENTROPY Procedure' is highlighted. The page content discusses the ENTROPY procedure, its features, and provides a link to the official SAS website (<http://www.sas.com>). The footer indicates the slide is from 'A. Holzinger 709.049' at '67/74'.

**Backup Slide: Example Tool for large data sets - Hadoop**

The slide shows a screenshot of the Apache Hadoop website. The header features the Hadoop logo. The main content area is titled 'Welcome to Apache™ Hadoop™!' and contains information about the project, including its history and a list of sub-project links like 'HDFS', 'MapReduce', 'HBase', 'Avro', 'Cassandra', 'Hive', 'Mahout', 'Pig', and 'ZooKeeper'. The footer indicates the slide is from 'A. Holzinger 709.049' at '68/74'.

**Backup Slide: Methods for Mining ..**

The diagram illustrates the hierarchy of mining methods. At the top is 'Topological Mining' (represented by a grid). Below it is 'Text Mining' (represented by a speech bubble). Further down is 'Data Mining' (represented by a cylinder). The bottom layer is divided into 'Standardized' (left) and 'Non-Standardized' (right). The left side is labeled 'Well-Structured' and the right side is labeled 'Weakly-Structured'. The bottom row is labeled 'reality (t<sub>1</sub>)' on the left and 'reality (t<sub>2</sub>)' on the right. The slide is attributed to Holzinger, A. (2011).

**Backup Slide: Excursion: How to get theories?**

The diagram compares two approaches to theory development. On the left, 'positivism' is shown as a process where 'theories and models' are abstracted from 'reality (t<sub>1</sub>)' and then concretized back into 'reality (t<sub>2</sub>)'. On the right, 'constructionism' is shown as a process where 'theories and models' are developed directly within 'reality (t<sub>2</sub>)'. The slide is attributed to Rauterberg, M. (2006).

**Backup Slide: The DIKW Model (1/4)**

A cartoon illustration showing four stages of the DIKW model. Stage 1: 'Data' (a caveman holding a spear). Stage 2: 'INFORMATION...' (the same caveman looking at a map). Stage 3: 'KNOWLEDGE...' (the caveman looking thoughtful). Stage 4: 'WISDOM!' (the caveman looking wise). The slide is attributed to Cleveland H. "Information as Resource", The Futurist, December 1982 p 34-39.

**Backup Slide: The DIKW Model (2/4)**

A diagram of the DIKW hierarchy represented as a triangle. The top vertex is 'Wisdom' (non-algorithmic, non-programmable). The bottom vertex is 'Data' (algorithmic, programmable). The middle vertex is 'Knowledge (actionable information)'. The sides of the triangle are labeled 'Information (data "in formation")'. The slide is attributed to Rowley, J. (2007).

