



Andreas Holzinger
Machine Learning & Knowledge Extraction
for Health Informatics
University of Verona
Module 2 - Day 2 – April 2017



Health Data Jungle: From the underlying physics of data to the Kullback-Leibler Divergence

a.holzinger@hci-kdd.org

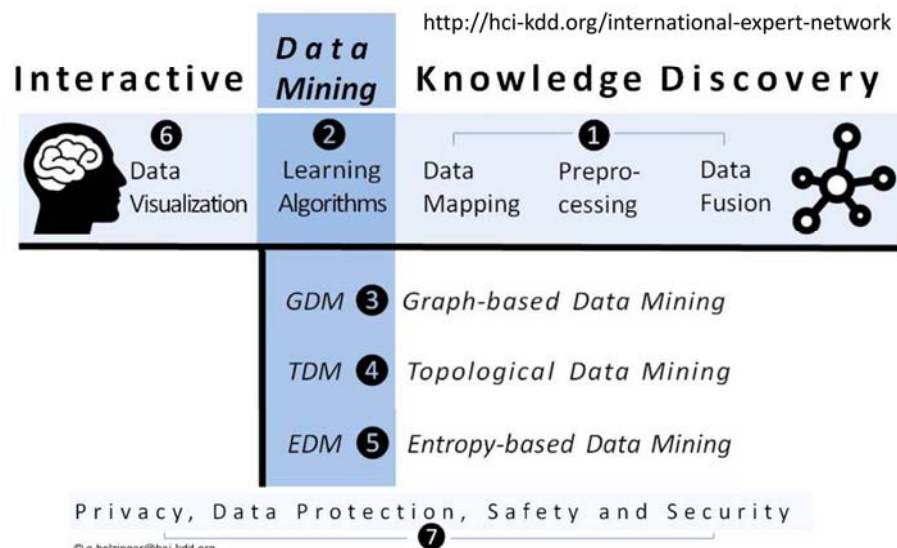
<http://hci-kdd.org/mini-make-machine-learning-knowledge-extraction-health>



Holzinger Group, hci-kdd.org

1

MAKE Health Verona 02



Holzinger, A. 2014. Trends in Interactive Knowledge Discovery for Personalized Medicine: Cognitive Science meets Machine Learning. IEEE Intelligent Informatics Bulletin, 15, (1), 6-14.

Holzinger Group, hci-kdd.org

3

MAKE Health Verona 02

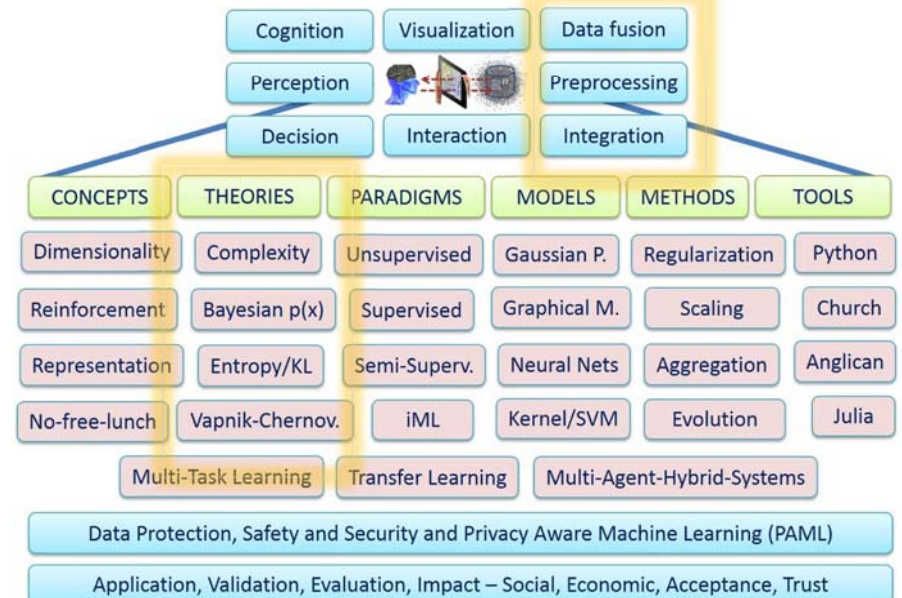
- 01 Data – the underlying physics of data
- 02 Biomedical data sources – taxonomy of data
- 03 Data integration, mapping, fusion
- 04 Probabilistic Information
- 05 Information Theory – Information Entropy
- 06 Cross- Entropy - Kullback-Leibler Divergence



Holzinger Group, hci-kdd.org

2

MAKE Health Verona 02



Holzinger, A. 2016. Machine Learning for Health Informatics. In: LNCS 9605, pp. 1-24, doi:10.1007/978-3-319-50478-0_1.

Holzinger Group, hci-kdd.org

4

MAKE Health Verona 02

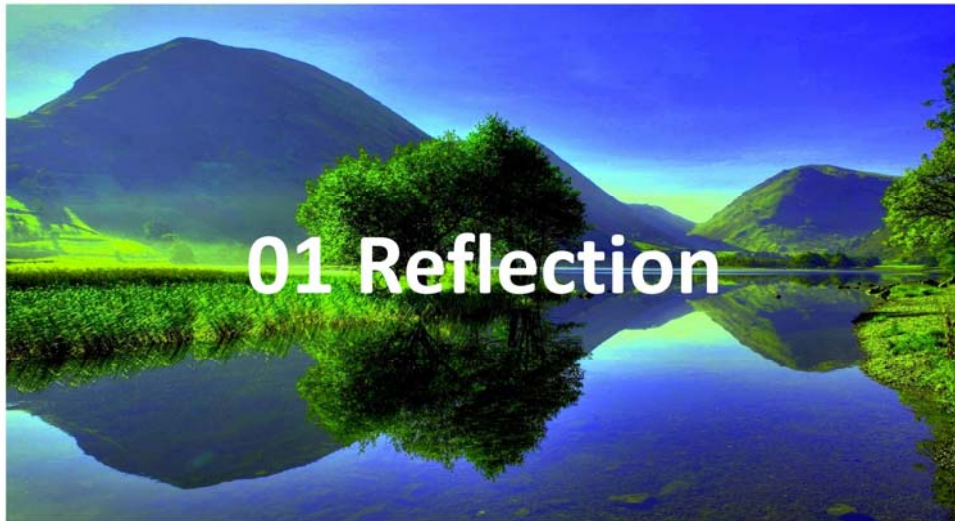


Image source: <http://www.hutui6.com/reflection-wallpapers.html>

Question: Where is the Biologist in this image?



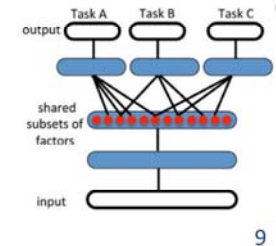
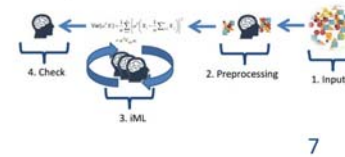
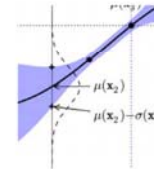
Image source: <http://www.efmc.info/medchemwatch-2014-1/lab.php>

Domingos, P. 2015. The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World, Penguin UK.

Warm-up Quiz



$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) * p(\theta)}{p(\mathcal{D})} \quad 3$$



Repetition of Bayes - on the work of Laplace

What is the simplest mathematical operation for us?

$$p(x) = \sum_y (p(x, y)) \quad (1)$$

How do we call repeated adding?

$$p(x, y) = p(y|x) * p(y) \quad (2)$$

Laplace (1773) showed that we can write:

$$p(x, y) * p(y) = p(y|x) * p(x) \quad (3)$$

Now we introduce a third, more complicated operation:

$$\frac{p(x, y) * p(y)}{p(y)} = \frac{p(y|x) * p(x)}{p(y)} \quad (4)$$

We can reduce this fraction by $p(y)$ and we receive what is called Bayes rule:

$$p(x, y) = \frac{p(y|x) * p(x)}{p(y)} \quad p(h|d) = \frac{p(d|h)p(h)}{p(d)} \quad (5)$$

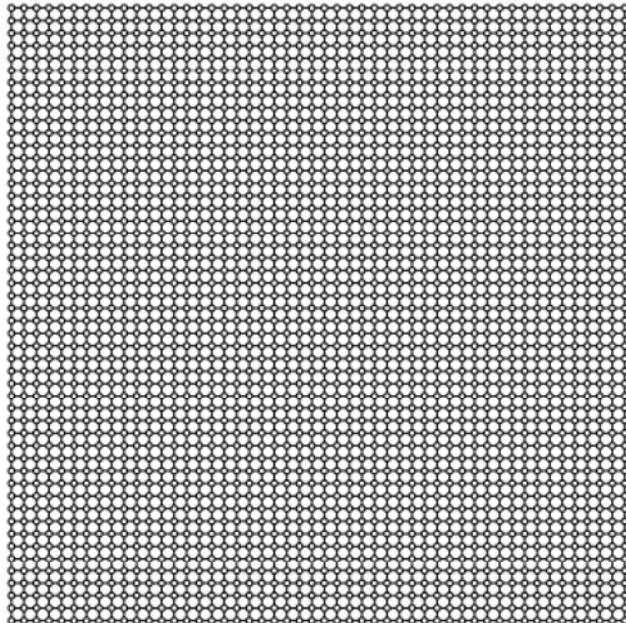
- Your MD has bad news and good news for you.
- Bad news first: You are tested positive for a serious disease, and the test is 99% accurate (T)
- Good news: It is a rare disease, striking 1 in 10,000 (D)
- **How worried would you now be?**



$$\text{posterior } p(x) = \frac{\text{likelihood} * \text{prior } p(x)}{\text{evidence}} \quad p(h|d) = \frac{p(d|h)p(h)}{p(d)}$$

$$p(T = 1|D = 1) = p(d|h) = 0,99 \text{ and } p(D = 1) = p(h) = 0,0001$$

$$p(D = 1 | T = 1) = \frac{(0,99) * (0,0001)}{(1-0,99) * (1-0,0001) + 0,99 * 0,0001} = 0,0098$$

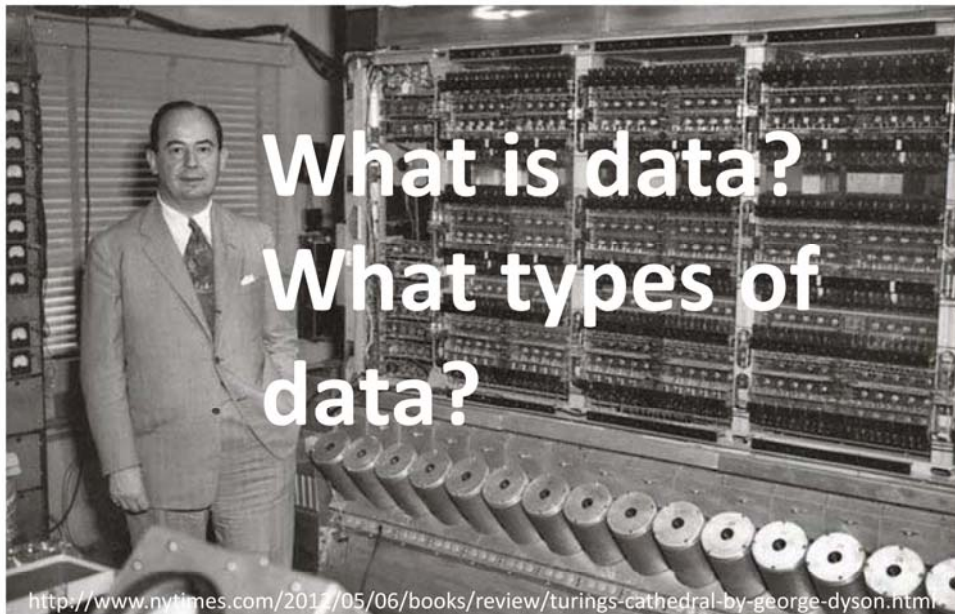


- Heterogeneous, distributed, inconsistent data sources (need for **data integration & fusion**) [1]
- **Complex data** (high-dimensionality – challenge of dimensionality reduction and visualization) [2]
- Noisy, uncertain, missing, dirty, and imprecise, imbalanced data (challenge of **pre-processing**)
- The discrepancy between data-information-knowledge (**various definitions**)
- **Big data** sets in high-dimensions (manual handling of the data is 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):11.
2. Hund, M., Sturm, W., Schreck, T., Ullrich, T., Keim, D., Majnarić, L. & Holzinger, A. 2015. Analysis of Patient Groups and Immunization Results 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.

**Statistical inference &
Decision support:
Better a good solution
in time,
than a perfect solution never ...**

01 The underlying physics of data



- 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 ($>> \mathbb{R}^{100}$)
- Problem: not noise, but complexity
- Much structure, but the structure can **not** be represented by a simple model

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

Diagnosis [E01]
Diagnostic Techniques and Procedures [E01.370]
Mass Screening [E01.370.500]

Anonymous Testing [E01.370.500.174]
Mass Chest X-Ray [E01.370.500.500]
Multiphasic Screening [E01.370.500.540]
▶ Neonatal Screening [E01.370.500.580]

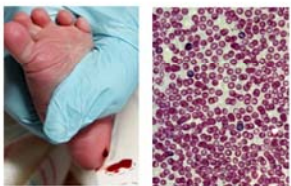
Diagnosis [E01]
Laboratory Techniques and Procedures [E01.450]

Age Determination by Skeleton [E01.450.074]
Clinical Chemistry Tests [E01.450.150] +
Cytodiagnosis [E01.450.230] +
Hematologic Tests [E01.450.375] +
Immunologic Tests [E01.450.495] +
Metabolic Clearance Rate [E01.450.520]
▶ Neonatal Screening [E01.450.560]
Occult Blood [E01.450.575]
Parasite Egg Count [E01.450.600]
Pregnancy Tests [E01.450.620] +
Radioligand Assay [E01.450.650]
Semen Analysis [E01.450.752] +
Sex Determination Analysis [E01.450.855]
Sex Determination by Skeleton [E01.450.860]
Specimen Handling [E01.450.865] +
Urinalysis [E01.450.890]

Medical Subject Headings MeSH

Newborn screening	
Intervention	
	
MeSH	D015997
MedlinePlus	007257

http://www.nlm.nih.gov/cgi/mesh/2011/MB_cgi?mode=&index=15177&view=expanded#TreeE01.370.500.580

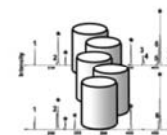


Amino acids (symbols)	Fatty acids (symbols)	Fatty acids (symbols)
Alanine (Ala)	Free carnitine (C0)	Hexadecenyl-carnitine (C16:1)
Arginine (Arg)	Acetyl-carnitine (C2)	Octadecenyl-carnitine (C18:1)
Argininosuccinate (Argsuc)	Propionyl-carnitine (C3)	Decenyl-carnitine (C10:2)
Citrulline (Cit)	Butyryl-carnitine (C4)	Tetradecenyl-carnitine (C14:2)
Glutamate (Glu)	Isovaleryl-carnitine (C5)	Octadecenyl-carnitine (C18:2)
Glycine (Gly)	Hexanoyl-carnitine (C6)	Hydroxy-isovaleryl-carnitine (C5-OH)
Methionine (Met)	Octanoyl-carnitine (C8)	Hydroxytetradecenyl-carnitine (C14-OH)
Ornithine (Orn)	Decanoyl-carnitine (C10)	Hydroxypalmitoyl-carnitine (C16-OH)
Phenylalanine (Phe)	Dodecanoyl-carnitine (C12)	Hydroxypalmitoyl-carnitine (C16:1-OH)
Pyroglutamate (PyrGlu)	Myristoyl-carnitine (C14)	Hydroxyl-carnitine (C18:1-OH)
Serine (Ser)	Hexadecenyl-carnitine (C16)	Decarboxyl-butyryl-carnitine (C4-DC)
Tyrosine (Tyr)	Octadecenyl-carnitine (C18)	Glutaryl-carnitine (C5-DC)
Valine (Val)	Tiglyl-carnitine (C3:1)	Methylglutaryl-carnitine (C6-DC)
Leucine + Isoleucine (Xci)	Decenyl-carnitine (C10:1)	Methylglutaryl-carnitine (C12-DC)
	Myristoyl-carnitine (C14:1)	

Fourteen amino acids and 29 fatty acids are analyzed from a single blood spot using MS/MS. The concentrations are given in $\mu\text{mol/L}$.



Yao, Y., Bowen, B. P., Baron, D. & Poznanski, D. 2015. SciDB for High-Performance Array-Structured Science Data at NERSC. Computing in Science & Engineering, 17, (3), 44-52, doi:10.1109/MCSE.2015.43.



DB of high-dimensional metabolic data including cases designated as PAHD (n=94), MCADD (n=63) and 3-MCCD (n=22), and a randomly sampled number of controls (n=1241)

Construction of classification models

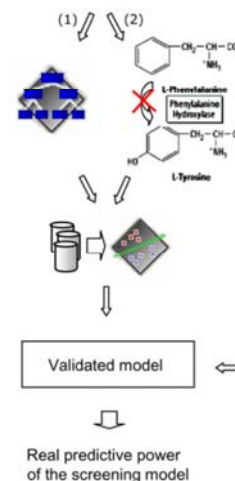
(1) decision tree paradigm with internal feature selection strategy

(2) Logistic regression analysis with expert knowledge (diagnostic flags) as model input variables

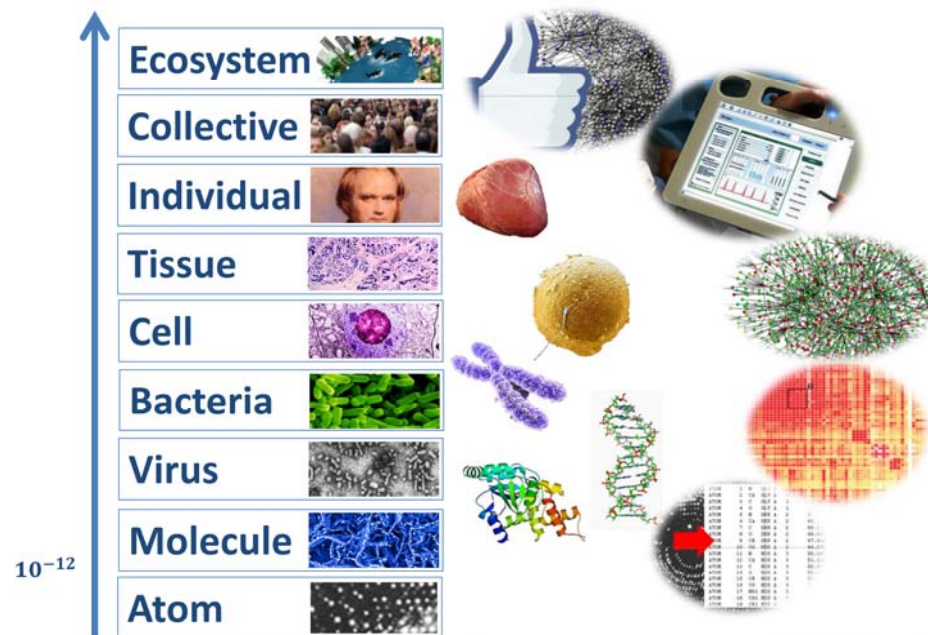
Training and 10-fold-cross validation

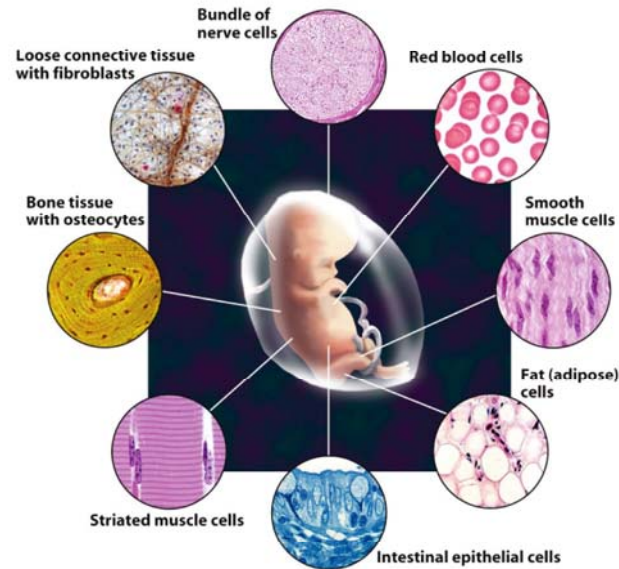
Larger database of control individuals (n=98,411) in order to estimate the specificity of a representative screening population

Baumgartner, C., Bohm, C. & Baumgartner, D. 2005. Modelling of classification rules on metabolic patterns including machine learning and expert knowledge. Journal of Biomedical Informatics, 38, (2), 89-98, doi:10.1016/j.jbi.2004.08.009.



02 Biomedical data sources: Taxonomy of data



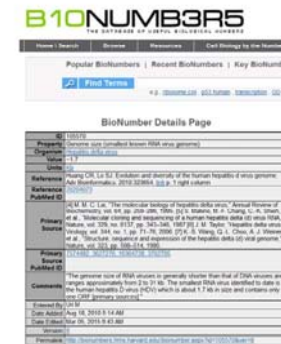


Karp, G. 2010. Cell and Molecular Biology: Concepts and Experiments, Gainesville, John Wiley.

Holzinger Group, hci-kdd.org

21

MAKE Health Verona 02



bionumbers.hms.harvard.edu/

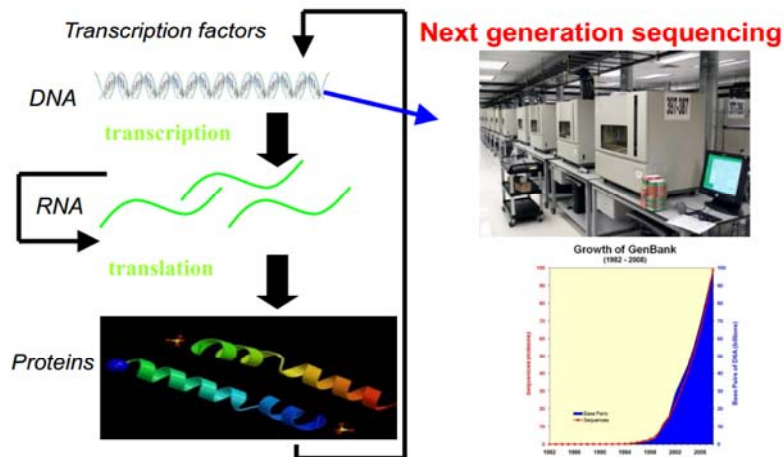
<http://book.bionumbers.org/how-many-genes-are-in-a-genome/>

Holzinger Group, hci-kdd.org

Organism	# of protein-coding genes	# of genes naive estimate: (genome size /1000)	BNID
HIV 1	9	10	105769
Influenza A virus	10-11	14	105767
Bacteriophage λ	66	49	105770
Epstein Barr virus	80	170	103246
Buchnera sp.	610	640	105757
T. maritima	1,900	1,900	105766
S. aureus	2,700	2,900	105500
V. cholerae	3,900	4,000	105760
B. subtilis	4,400	4,200	111448
E. coli	4,300	4,600	105443
S. cerevisiae	6,600	12,000	105444
C. elegans	20,000	100,000	101364
A. thaliana	27,000	140,000	111380
D. melanogaster	14,000	140,000	111379
F. rubripes	19,000	400,000	111375
Z. mays	33,000	2,300,000	110565
M. musculus	20,000	2,800,000	100308
H. sapiens	21,000	3,200,000	100399, 111378
T. aestivum (hexaploid)	95,000	16,800,000	105448, 102713

22

MAKE Health Verona 02

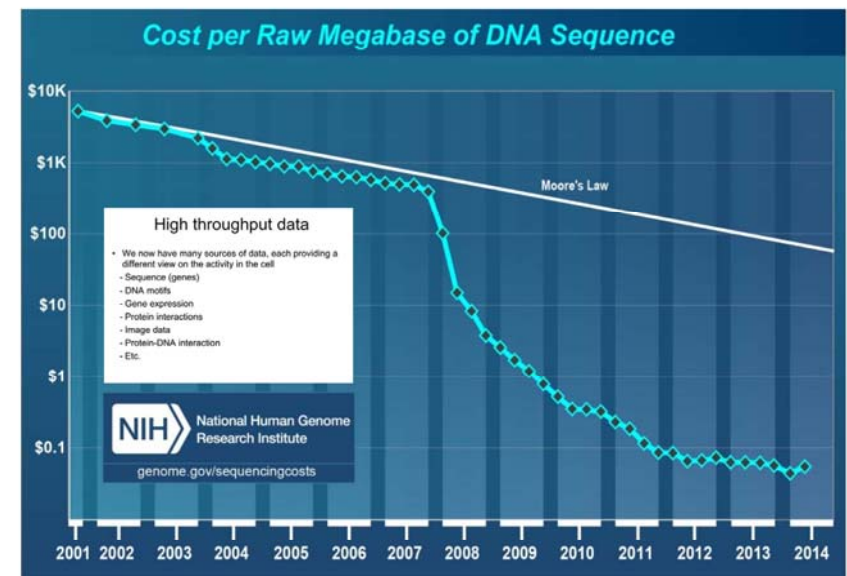


Navlakha, S. & Bar-Joseph, Z. 2011. Algorithms in nature: the convergence of systems biology and computational thinking. *Molecular Systems Biology*, 7.

Holzinger Group, hci-kdd.org

23

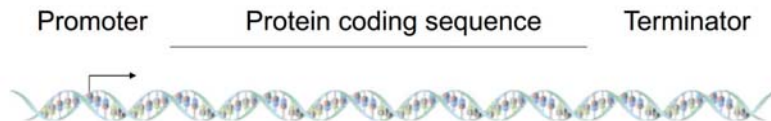
MAKE Health Verona 02



Holzinger Group, hci-kdd.org

24

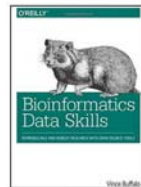
MAKE Health Verona 02



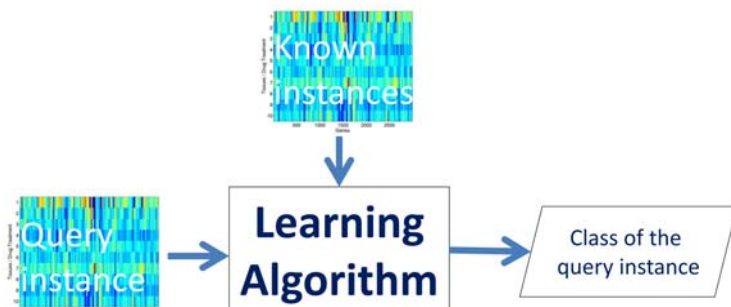
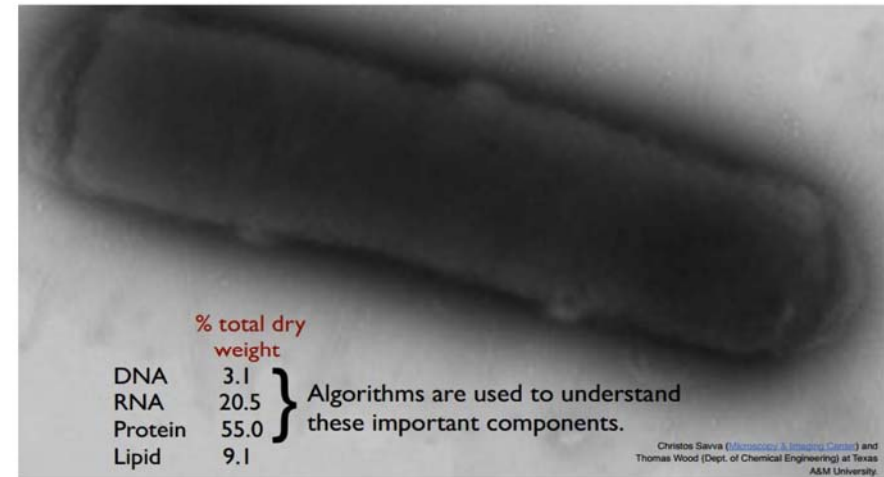
ATGAAGCTACTGTCTTCTATCGAACAAGCATCGCATATTTGCCGACTTAAAGGCTCAAG
 TGCTCCAAAGAAAAACCGAAGTGCCTCAAGTGTCTGAAGAACAAGTGGAGTGTGCTGCTAC
 TCTCCAAAGAAAAACCGAAGTGCCTCAAGTGTCTGAAGAACAAGTGGAGTGTGCTGCTAC
 CTGAAAGAGCTGGAACAGCTATTTCTACTGATTTTCTCGAAGAAGACCTTGACATGATT
 TTGAAATGGATTCTTACAGGATATAAAGCATTGTTAACAGGATATTGTACAAAGAT
 AATGTGAATAAAGATCGCTCAGCAGATAGATTGGCTTCAAGTGAAGACTGATATGCTCTA
 ACATTGAGACAGCATAGAATAAGTGCACATCATCGGAAGAGAGTAGTAACAAAGGT
 CAAAGACAGTTGACTGTATCGATTGACTCGGCAGCTCATGATGAAGTCCACAATCCG
 TTGGATTTTATGCCAGGGATGCTTCTCATGGATTGATTGGTCTGAAGAGGATGACATG
 TCGGATGGCTTGGCCTTCTGAAACGGACCCCAACAATAATGGTCTTTTGGCGACGGT
 TCTCTCTATGATTCTTCTGATCATTTGGCTTAAACCGGAAATACAGCACTCTAAC
 GTTAACAGGCTCCCGACCATGATTACGGGATAGATACAGCTTGGCTTCTGATCGCAACA
 TCCCGTTACTCAAAATATCTCAATAATTTTACCCTACTGCCCTATCGTGCACTCA
 CCGACGCTAATGATGTTGATAATAACCAAGATTGAAATCGGCTCGAAGGATCAATGGCAA
 ATCCTTTTAACTGCATATAGCCATTGGAGCCTGGGTATAGAGGGGGGAATCTACTGAT
 ATAGATGTTTTTACTATCAAAATGCTAAATCTCATTTGACGAGCAAGGTTCTCGAGTCA

		Second Letter				
		U	C	A	G	
1st letter	U	UUU Phe	UUC Ser	UAU Tyr	UGU Cys	U
	U	UUC Ser	UCC Ser	UAC Tyr	UGC Cys	C
	U	UUA Leu	UUA Leu	UAA Stop	UGA Stop	A
C	C	CUU Leu	CCU Pro	CAU His	CGU Arg	C
	C	CUC Leu	CCC Pro	CAC His	CGC Arg	G
	C	CUG Leu	CCG Pro	CAG Gln	CGA Arg	A
A	A	AUU Ile	ACU Thr	AAU Asn	AGU Ser	A
	A	AUC Ile	ACC Thr	AAC Asn	AGC Ser	C
	A	AUA Ile	ACA Thr	AAA Lys	AGA Arg	A
G	G	GUU Val	GCU Ala	GAU Asp	GGU Gly	G
	G	GUC Val	GCC Ala	GAC Asp	GGC Gly	A
	G	GUA Val	GCG Ala	GAG Asp	GGG Gly	C

For further reading this is recommended:
 Buffalo, V. 2015. Bioinformatics Data Skills:
 Reproducible and Robust Research with Open
 Source Tools, Sebastopol (CA), O'Reilly.



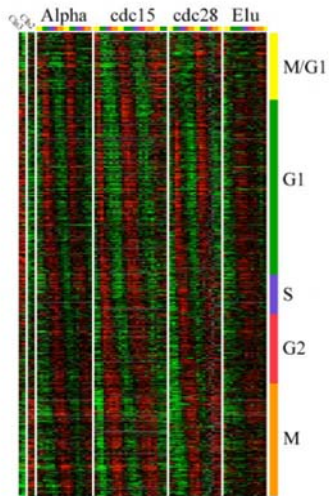
MAKE Health Verona 02



Features are key to learning
and understanding!

- Billions of biological data sets are openly available, here only some examples:
- General Repositories:
 - GenBank, EMBL, HMCA, ...
- Specialized by data types:
 - UniProt/SwissProt, MMMP, KEGG, PDB, ...
- Specialized by organism:
 - WormBase, FlyBase, NeuroMorpho, ...
- Details: <http://hci-kdd.org/open-data-sets>





- this figure depicts one yeast gene-expression data set
- each row represents a gene
- each column represents a measurement of gene expression (mRNA abundance) at some time point
- red indicates that a gene is being expressed more than some baseline; green means less

Figure from Spellman et al., Molecular Biology of the Cell, 9:3273-3297, 1998

- 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, numerical measurements (physiological data, lab results, vital signs, ...), recorded signals (ECG, EEG, ...), Images (x-ray, MR, CT, PET, ...) ; -omics

TU Taxonomy of data at Hospital Level

Clinical workplace data sources

- Med.docs: text (non-standardized (free-text), semi-structured, standard terminologies (ICD, SNOMED-CT))
- Measurements: lab results, ECG, EEG, EOG, ...
- Surveys, Clinical studies, trials

Image data sources

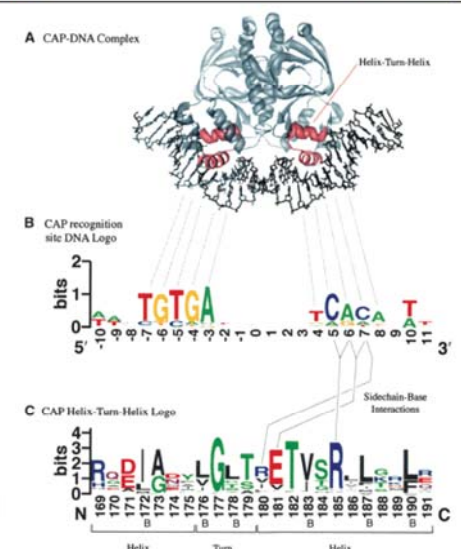
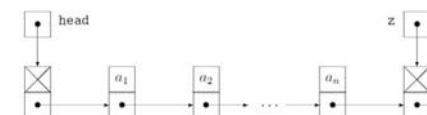
- Radiology: MRI (256x256, 200 slices, 16 bit per pixel, uncompressed, ~26 MB); CT (512x512, 60 slices, 16 bit per pixel, uncompressed ~32MB; MR, US;
- Digital Microscopy : WSI (15mm slide, 20x magn., 24 bits per pixel, uncompressed, 2,5 GB, WSI 10 GB; confocal laser scanning, etc.

-omics data sources

- Sanger sequencing, NGS whole genome sequencing (3 billion reads, read length of 36) ~ 200 GB; NGS exome sequencing ("only" 110,000,000 reads, read length of 75) ~7GB; Microarray, mass-spectrometry, gas chromatography, ...

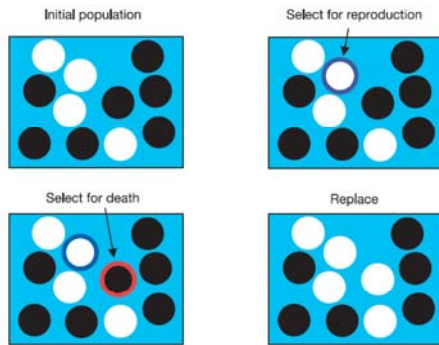
TU Example Data Structures (1/3): List

<pre> TYPE link = REF node ; node = RECORD key : ItemType; next : link; END; </pre>	<pre> key next </pre>	<pre> class link { ItemType key; link next; } </pre>
VAR p, q : link ;	p q	link p,q;
p := NEW(link);	p	p:=new link();
p->key:=x;	p	p.key:=x;
q := NEW(link) ;	q	q:=new link();

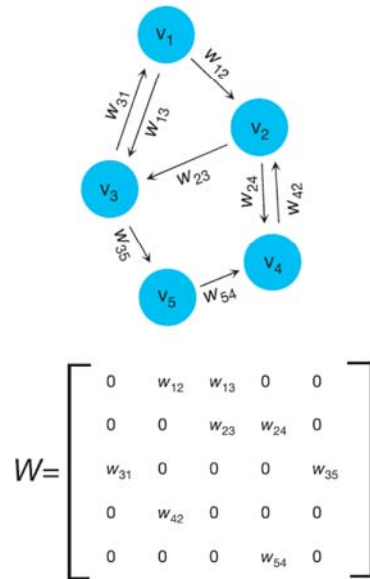


Crooks, G. E., Hon, G., Chandonia, J. M. & Brenner, S. E. (2004) WebLogo: A sequence logo generator. Genome Research, 14, 6, 1188-1190.

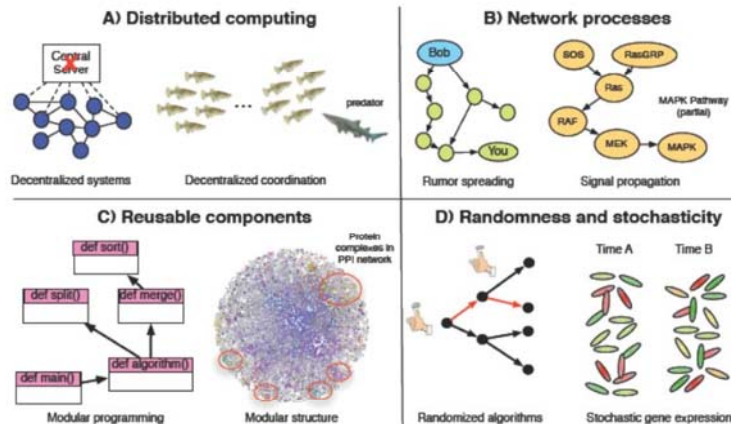
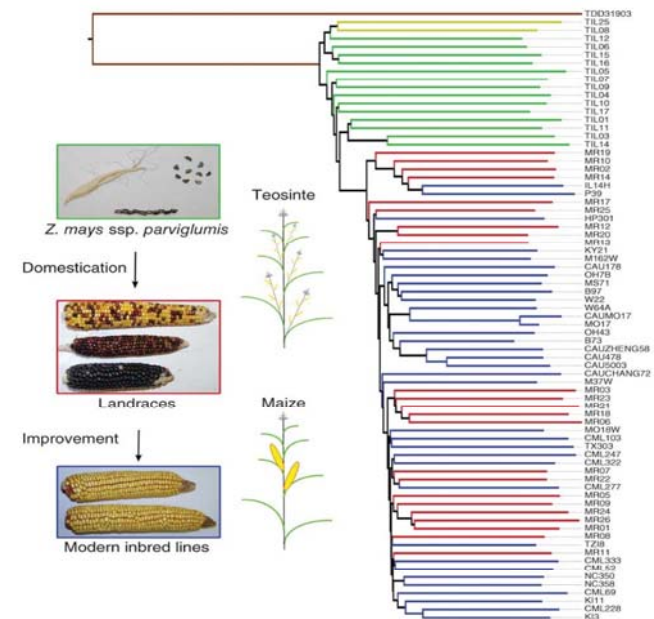
Evolutionary dynamics act on populations.
Neither genes, nor cells, nor individuals evolve;
only populations evolve.



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



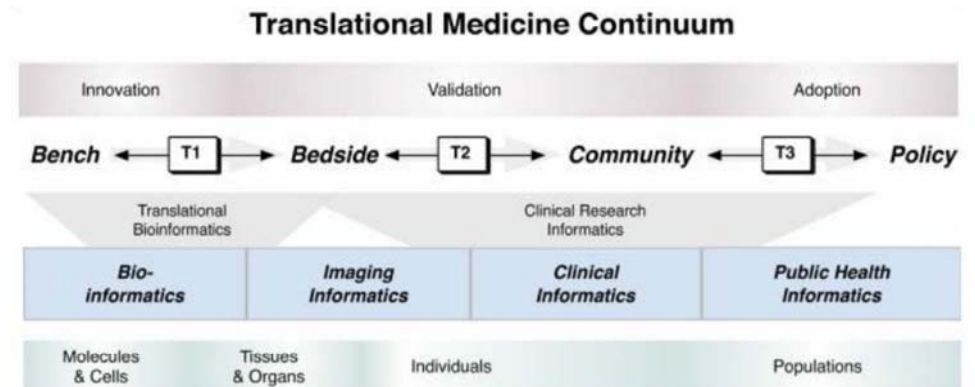
Hufford et. al.
2012. Comparative
population
genomics of maize
domestication and
improvement.
Nature Genetics,
44, (7), 808-811.



<http://cacm.acm.org/magazines/2015/1/181614-distributed-information-processing-in-biological-and-computational-systems/abstract>

Navlakha, S. & Bar-Joseph, Z. 2014. Distributed information processing in biological and computational systems. *Commun. ACM*, 58, (1), 94-102.

<https://www.youtube.com/watch?v=4u47nWHzqI4&feature=youtu.be>



Biomedical Informatics Continuum

Sarkar, I. 2010. Biomedical informatics and translational medicine. *Journal of Translational Medicine*, 8, (1), 2-12.

- Grand Challenges in this area:
 - – Production of Open Data Sets
 - – Synthetic data sets for learning algorithm testing
 - – Privacy preserving machine learning
 - – Data leak detection
 - – Data citation
 - – Differential privacy
 - – Anonymization and pseudonymization
 - – Evaluation and benchmarking

Please visit:

<http://hci-kdd.org/privacy-aware-machine-learning-for-data-science/>



03 Data Integration, mapping, fusion

Unsolved Problem: Data Integration and Data Fusion in the Life Sciences

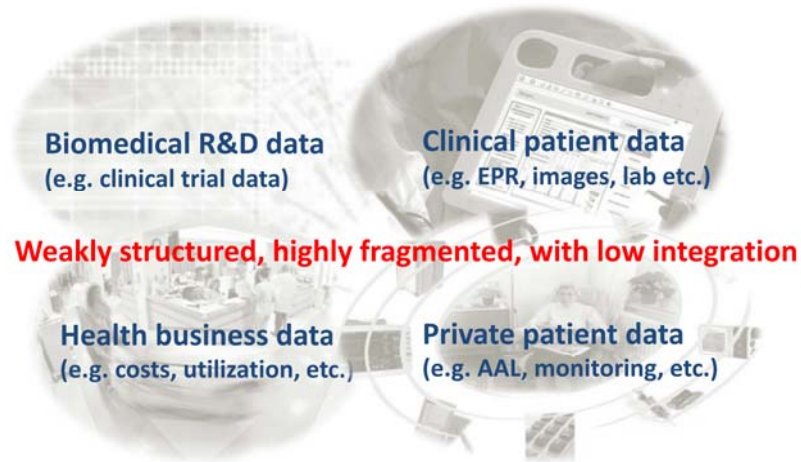
How to combine these different data types together to obtain a unified view of the activity in the cell is one of the major challenges of systems biology

Navlakha, S. & Bar-Joseph, Z. 2014. Distributed information processing in biological and computational systems. *Commun. ACM*, 58, (1), 94-102, doi:10.1145/2678280.

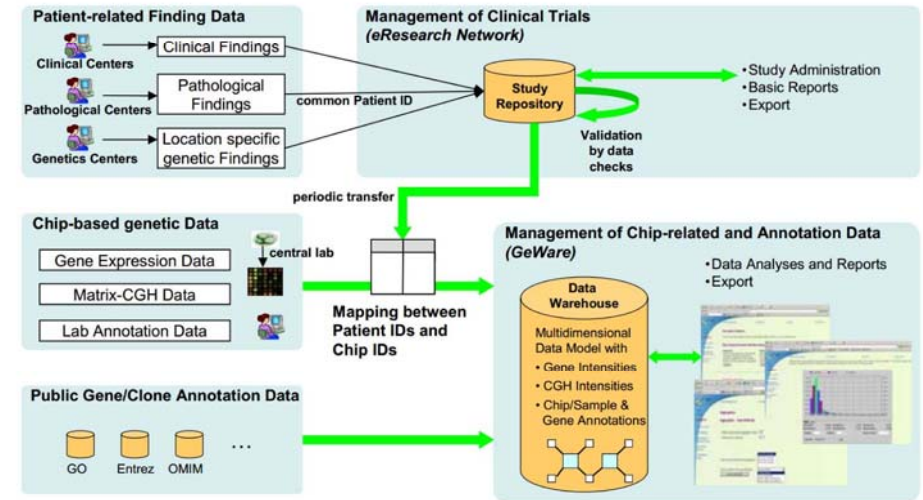


Our central hypothesis: Information may bridge this gap

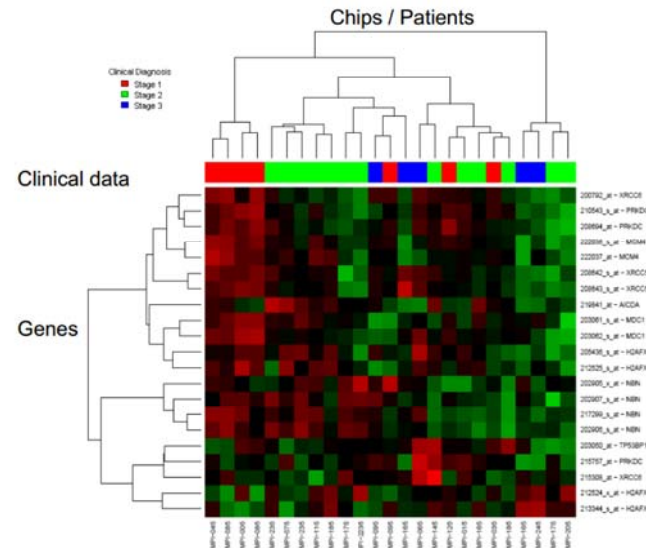
Holzinger, A. & Simon, K.-M. (eds.) 2011. *Information Quality in e-Health. Lecture Notes in Computer Science LNCS 7058, Heidelberg, Berlin, New York: Springer.*



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.

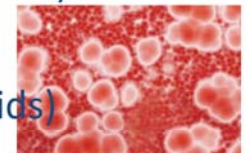


Kirsten, T., Lange, J. & Rahm, E. 2006. An integrated platform for analyzing molecular-biological data within clinical studies. *Current Trends in Database Technology—EDBT 2006*. Heidelberg: Springer, pp. 399-410, doi:10.1007/11896548_31.



Kirsten, T., Lange, J. & Rahm, E. 2006. An integrated platform for analyzing molecular-biological data within clinical studies. *Current Trends in Database Technology—EDBT 2006*. Heidelberg: Springer, pp. 399-410, doi:10.1007/11896548_31.

- 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)

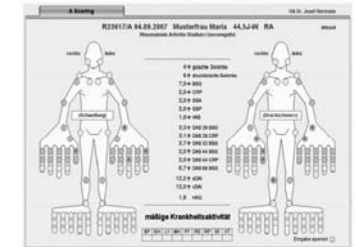
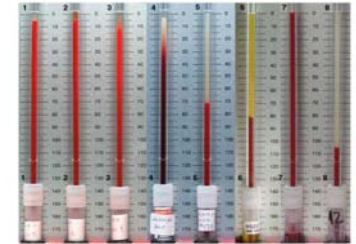


Genomics	Transcriptomics	Proteomics	Metabolomics	Protein-DNA interactions	Protein-protein interactions	Fluxomics	Phenomics
Genomics (sequence annotation)	• ORF validation • Regulatory element identification ¹⁴	• SNP effect on protein activity or abundance	• Enzyme annotation	• Binding-site identification ¹⁵	• Functional annotation ¹⁶	• Functional annotation	• Functional annotation ^{17,18,19} • Biomarkers ²⁰
	Transcriptomics (microarray, SAGE)	• Protein: transcript correlation ²¹	• Enzyme annotation ¹⁰⁰	• Gene-regulatory networks ²⁶	• Functional annotation ¹⁰¹ • Protein complex identification ¹⁰²		• Functional annotation ¹⁰²
		Proteomics (abundance, post-translational modification)	• Enzyme annotation ²⁰	• Regulatory complex identification	• Differential complex formation	• Enzyme capacity	• Functional annotation
			Metabolomics (metabolite abundance)	• Metabolic-transcriptional response		• Metabolic pathway bottlenecks	• Metabolic flexibility • Metabolic engineering ¹⁰³
				Protein-DNA interactions (ChIP-chip)	• Signalling cascades ^{101,102}		• Dynamic network responses ¹⁴
				Protein-protein interactions (yeast 2H, coAP-MS)			• Pathway identification activity ¹⁰⁴
						Fluxomics (isotopic tracing)	• Metabolic engineering
							Phenomics (phenotype arrays, RNAi screens, synthetic lethals)

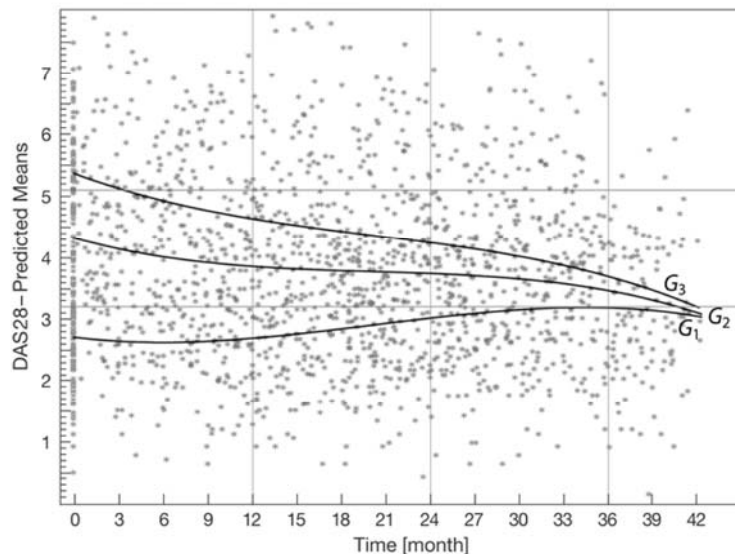


Joyce, A. R. & Palsson, B. Ø. 2006. The model organism as a system: integrating 'omics' data sets. *Nature Reviews Molecular Cell Biology*, 7, 198-210.

- 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.

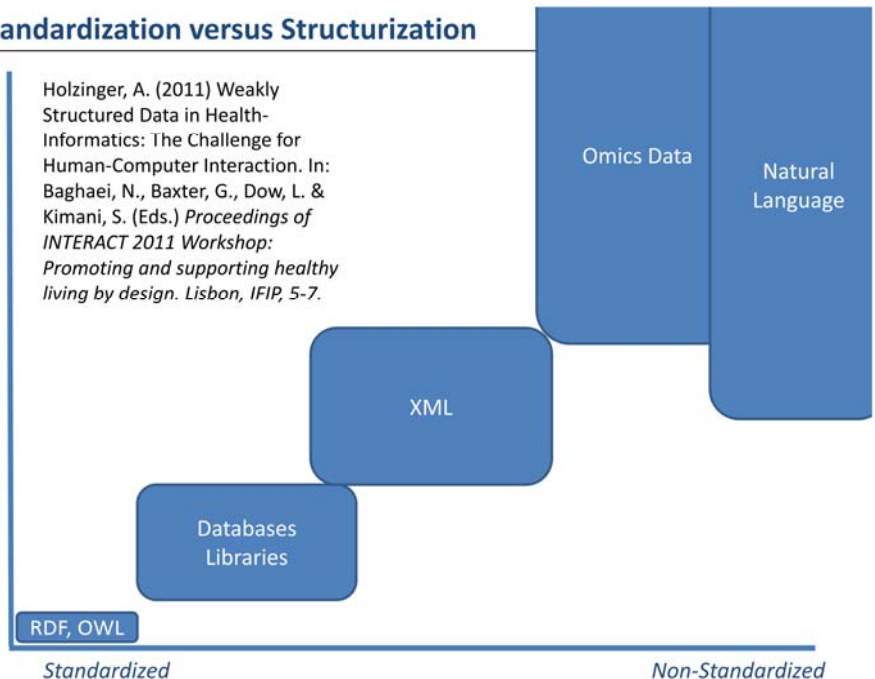


Simonic, K. M., Holzinger, A., Bloice, M. & Hermann, J. (2011). *Optimizing Long-Term Treatment of Rheumatoid Arthritis with Systematic Documentation*. *Pervasive Health - 5th International Conference on Pervasive Computing Technologies for Healthcare, Dublin, IEEE*, 550-554.

Weakly-Structured

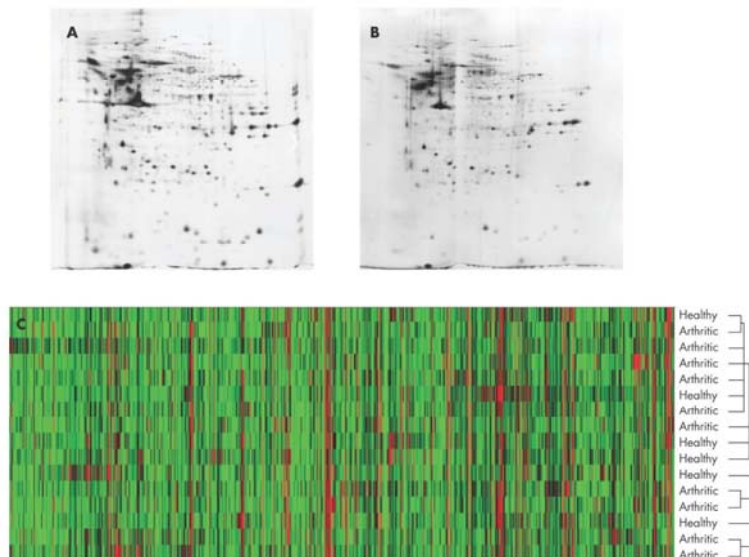
Well-Structured

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



- 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

TU 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.

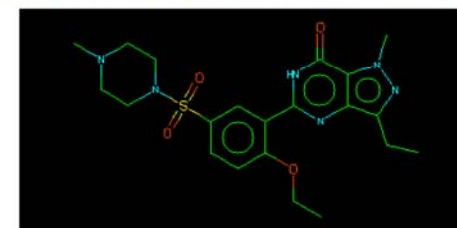
TU Example: 1-D data (univariate sequential data objects)

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

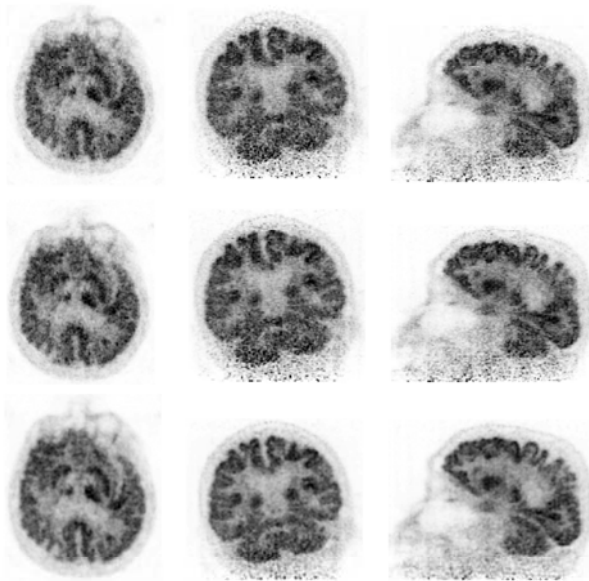
...is Comprehensive

...is Well Documented

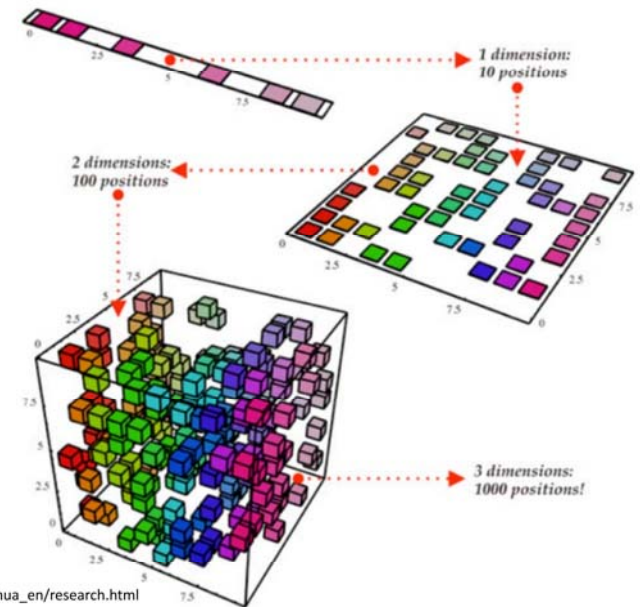
http://www.daylight.com/dayhtml_tutorials/languages/smiles/index.html

TU Example: 2.5-D data (structural information & metadata)

<http://www.pdb.org>

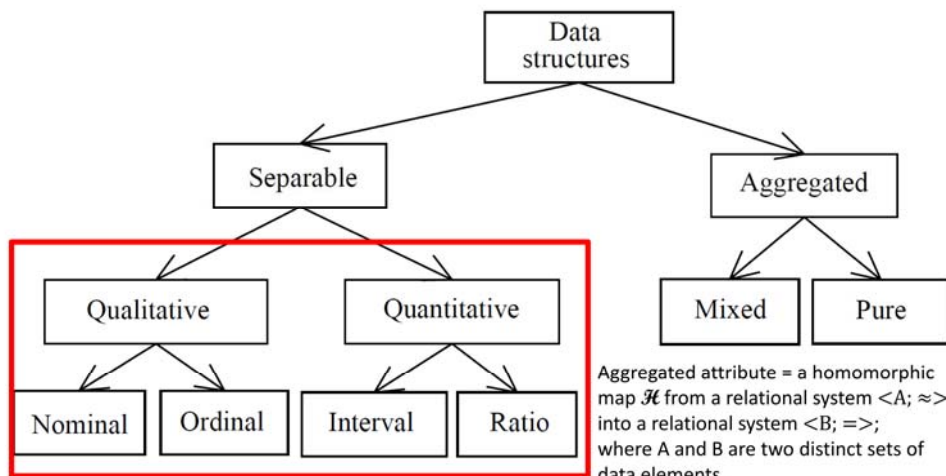


Scheins, J. J., Herzog, H. & Shah, N. J. (2011) Fully-3D PET Image Reconstruction Using Scanner-Independent, Adaptive Projection Data and Highly Rotation-Symmetric Voxel Assemblies. *Medical Imaging, IEEE Transactions on*, 30, 3, 879-892.



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



Dastani, M. (2002) The Role of Visual Perception in Data Visualization. *Journal of Visual Languages and Computing*, 13, 601-622.

Aggregated attribute = a homomorphic map \mathcal{H} from a relational system $\langle A; \approx \rangle$ into a relational system $\langle B; = \rangle$; where A and B are two distinct sets of data elements. This is in contrast with other attributes since the set B is the set of data elements instead of atomic values.

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 $x' = f(x)$ $x \dots \text{mono-tonic incr.}$	$x \mapsto f(x)$	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.	$=, \neq, >, <, -, +$
RATIO	Determination of equality or ratios	Similarity $x' = ax$	$x \mapsto rx$	Coefficient of variation	$=, \neq, >, <, -, +, *, \div$

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

- Bridging the gap between natural sciences and clinical medicine (who has seen genomics and patient data integrated in routine?)
- Organizational barriers, data provenance, data ownership, privacy, accessibility, usability, fair use of data, security, safety, data protection
- Combine Ontologies with Machine Learning
- Stochastic Ontologies, Ontology learning
- Integration of data from wet-labs with in-silico experimental data (e.g. tumor growth simulation)

- Boolean models
- Algebraic models
- Probabilistic models *)

*) Our probabilistic models describes data which we can observe from our environment – and if we use the mathematics of probability theory, in order to express the uncertainties around our model then the inverse probability allows us to infer unknown unknowns ... learning from data and making predictions – the core essence of machine learning and of vital importance for health informatics

Ghahramani, Z. 2015. Probabilistic machine learning and artificial intelligence. *Nature*, 521, (7553), 452-459, doi:10.1038/nature14541.

04 Probabilistic Information $p(x)$



Lane, N. & Martin, W. (2010) The energetics of genome complexity. *Nature*, 467, 7318, 929-934.

- **Communication** (Hartley, Nyquist, Shannon)
- **Coding Theory** (Fano, Hamming, Reed, Solomon)
- **Cryptography** (Hellman, Rivest, Shamir, Adleman)
- **Complexity** (Kolmogorov, Chaitin) **Computation, Chaos**
- **Cybernetics** (Wiener, von Neumann, Langton)
- **Foundations** (Brillouin, Bennet, Landauer)
- **Canonical Quantum Gravity** (Wheeler, De-Witt)
- **Metabiology** (Conrad, Chaitin)

Unification via Information (Carlo Rovelli's books)

Universe's ultimate mechanism for existence might be Information: "it from bit" (Wheeler's last speculation)

Manca, V. 2013. Infobiotics: Information in Biotic Systems, Heidelberg, Springer, doi:10.1007/978-3-642-36223-1.

Probabilistic Information $p(x)$



Bayes, T. (1763). An Essay towards solving a Problem in the Doctrine of Chances (Postum communicated by Richard Price). Philosophical Transactions, 53, 370-418.

$$p(x_i) = \sum P(x_i, y_j) \quad \text{Thomas Bayes 1701 - 1761} \quad p(x_i, y_j) = p(y_j|x_i)P(x_i)$$

Bayes' Rule is a corollary of the Sum Rule and Product Rule:

$$p(x_i|y_j) = \frac{p(y_j|x_i)p(x_i)}{\sum p(x_i, y_j)p(x_i)}$$

Barnard, G. A., & Bayes, T. (1958). Studies in the history of probability and statistics: IX. Thomas Bayes's essay towards solving a problem in the doctrine of chances. Biometrika, 45(3/4), 293-315.

Bayes Law of Total Probability = data modelling

Bayes' Rule in words

d ... data; h ... hypothesis

$H = \{H_1, H_2, \dots, H_n\}$... Hypothesis space

$\forall h, d \dots$

$$p(h|d) = \frac{p(d|h)p(h)}{\sum_{h \in H} p(d|h')p(h')}$$

Likelihood

Prior Probability

Posterior Probability

Evidence = marginal likelihood

Sum over space of alternative hypotheses

Always remember:

The inverse probability allows to infer unknowns, learn from data and make **predictions**:

1) Maximum-Likelihood Learning

finds a parameter setting, that maximizes the $p(x)$ of the data: $P(\mathcal{D}|\theta)$

2) Maximum a Posteriori Learning (e.g. for MCMC)

assumes a prior over the model parameters $P(\theta)$ and finds a parameter setting that maximizes the posterior: $P(\theta|\mathcal{D}) \propto P(\theta)P(\mathcal{D}|\theta)$

3) Bayesian Learning

assumes a prior over the model parameters and computes the posterior distribution $P(\theta|\mathcal{D})$

- General setting:
 - Given a (hypothesized & probabilistic) model that governs the random experiment
 - The model gives a probability of any data $p(D| \theta)$ that depends on the parameter θ
 - Now, given actual sample data $X = \{x_1, \dots, x_n\}$, what can we say about the value of θ ?
- Intuitively, take your best guess of θ
- “best” means “best explaining/fitting the data”
- Generally an optimization problem

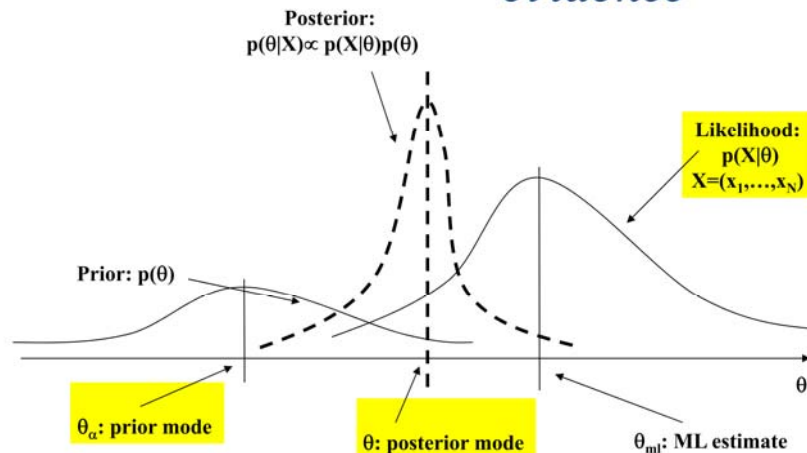
- 1) Maximum likelihood estimation (given X)
 - “Best” means “data likelihood reaches maximum”

$$\hat{\theta} = \arg \max_{\theta} P(X|\theta)$$
 - Problem: massive amount of data necessary**
- 2) Bayesian estimation (use posterior)

$$\hat{\theta} = \arg \max_{\theta} P(X|\theta) = \arg \max_{\theta} P(X|\theta) P(\theta)$$
 - “Best” means being consistent with our “prior” knowledge and explaining data well
 - Problem: how to define prior?**

An example can be found in: Banerjee, O., El Ghaoui, L. & D'aspremont, A. 2008. Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data. *The Journal of Machine Learning Research*, 9, 485-516. Available via: <http://arxiv.org/pdf/0707.0704>

$$\text{posterior } p(\theta) = \frac{\text{likelihood} * \text{prior } p(\theta)}{\text{evidence}}$$



For more basic information: Bishop, C. M. 2007. *Pattern Recognition and Machine Learning*, Springer.
For application examples in Text processing refer to: Jiang, J. & Zhai, C. X. 2007. An empirical study of tokenization strategies for biomedical information retrieval. *Information Retrieval*, 10, (4-5), 341-363.

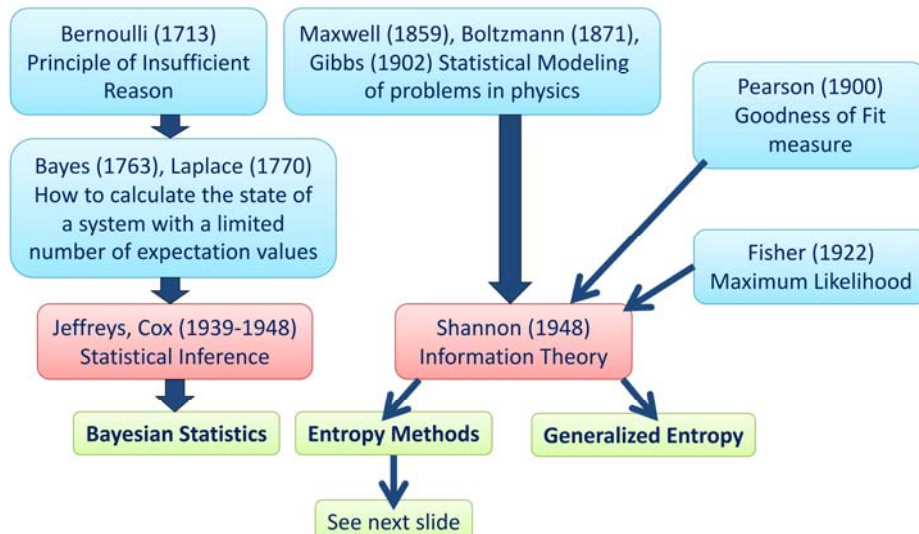
05 Information Theory & Entropy

- Information is the reduction of uncertainty
- If something is 100 % certain its uncertainty = 0
- Uncertainty is a max. if all choices are equally probable
- Uncertainty (as information) sums up for independent sources



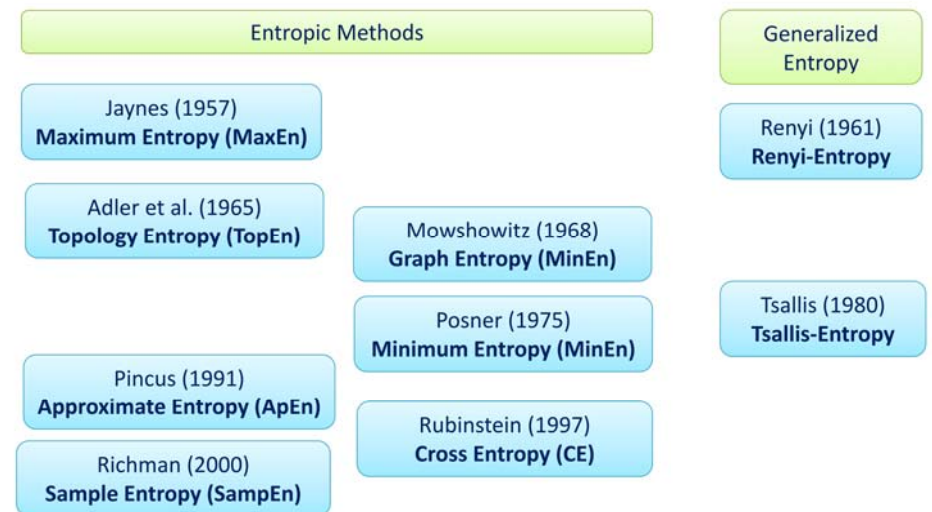
<http://www.scottaaronson.com>

TU An overview on the History of Entropy

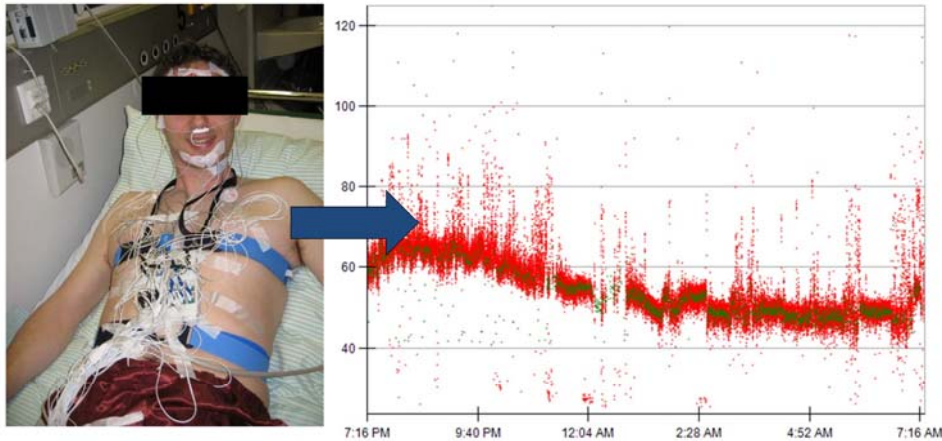


confer also with: Golan, A. (2008) Information and Entropy Econometric: A Review and Synthesis. *Foundations and Trends in Econometrics*, 2, 1-2, 1-145.

TU Towards a Taxonomy of Entropic Methods

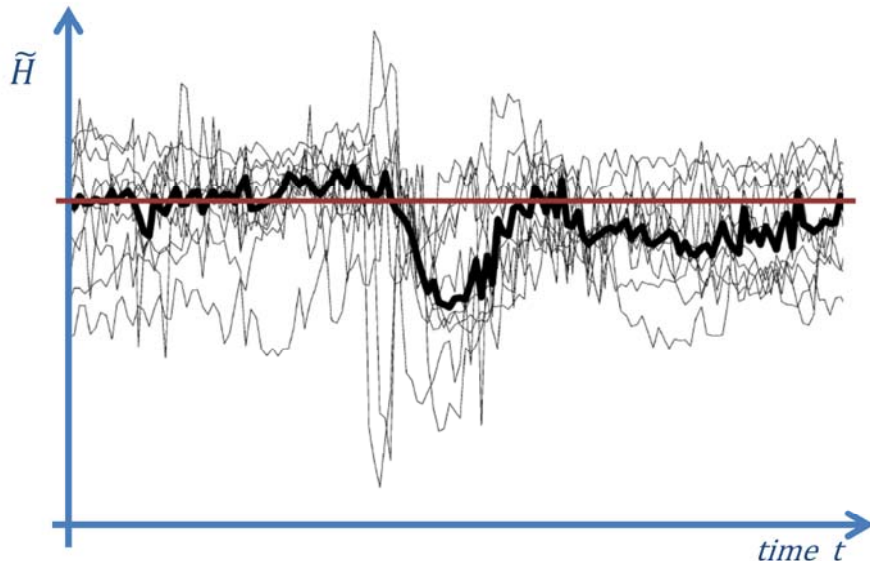


Holzinger, A., Hörtenhuber, M., Mayer, C., Bachler, M., Wassertheurer, S., Pinho, A. & Koslicki, D. 2014. On Entropy-Based Data Mining. In: Holzinger, A. & Jurisica, I. (eds.) *Lecture Notes in Computer Science, LNCS 8401*. Berlin Heidelberg: Springer, pp. 209-226.



Holzinger, A., Stocker, C., Bruschi, M., Auinger, A., Silva, H., Gamboa, H. & Fred, A. 2012. On Applying Approximate Entropy to ECG Signals for Knowledge Discovery on the Example of Big Sensor Data. In: Huang, R., Ghorbani, A., Pasi, G., Yamaguchi, T., Yen, N. & Jin, B. (eds.) *Active Media Technology, Lecture Notes in Computer Science, LNCS 7669. Berlin Heidelberg: Springer*, pp. 646-657.

EU Project EMERGE (2007-2010)



$$\text{Let: } \langle x_n \rangle = \{x_1, x_2, \dots, x_N\}$$

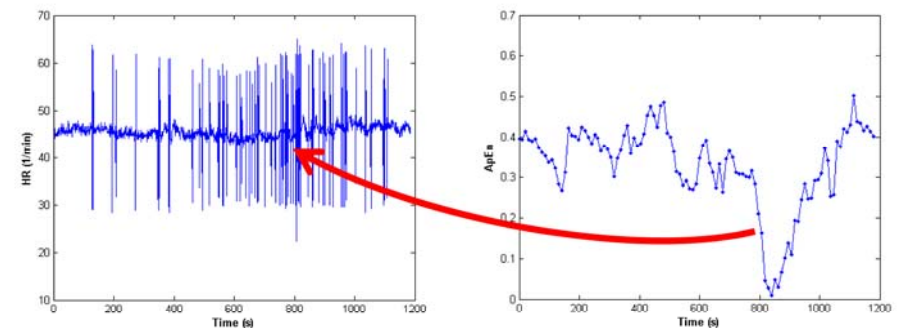
$$\vec{X}_i = (x_i, x_{(i+1)}, \dots, x_{(i+m-1)})$$

$$\|\vec{X}_i, \vec{X}_j\| = \max_{k=1,2,\dots,m} (|x_{(i+k-1)} - x_{(j+k-1)}|)$$

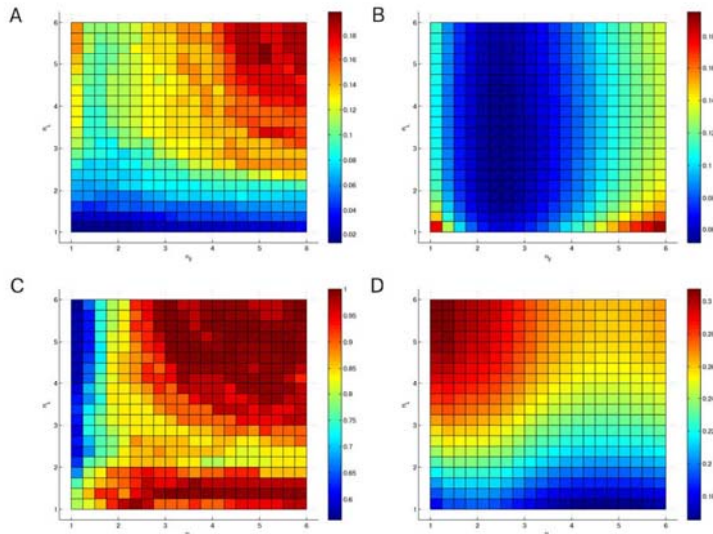
$$\tilde{H}(m, r) = \lim_{N \rightarrow \infty} [\phi^m(r) - \phi^{m+1}(r)]$$

$$C_r^m(i) = \frac{N^m(i)}{N - m + 1} \quad \phi^m(r) = \frac{1}{N - m + 1} \sum_{t=1}^{N-m+1} \ln C_r^m(i)$$

Pincus, S. M. (1991) Approximate Entropy as a measure of system complexity. *Proceedings of the National Academy of Sciences of the United States of America*, 88, 6, 2297-2301.



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

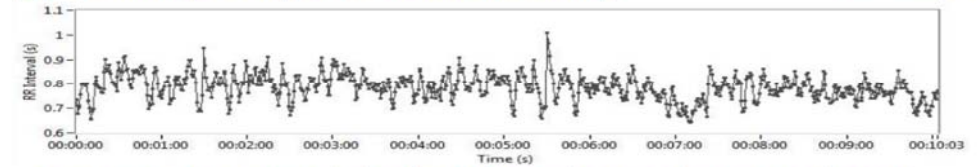


Mayer, C., Bachler, M., Hortenhuber, M., Stocker, C., Holzinger, A. & Wassertheurer, S. 2014. Selection of entropy-measure parameters for knowledge discovery in heart rate variability data. *BMC Bioinformatics*, 15, (Suppl 6), S2, doi:doi:10.1186/1471-2105-15-S6-S2.

Holzinger Group, hci-kdd.org

77

MAKE Health Verona 02



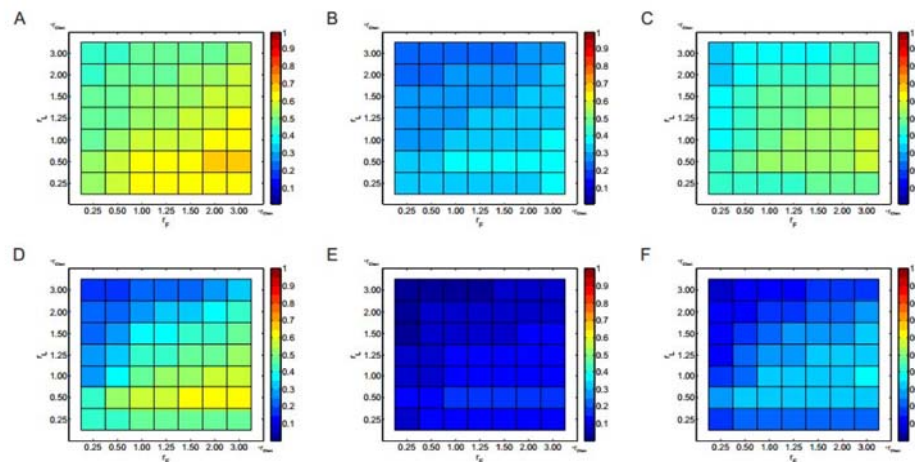
- Heart Rate Variability (HRV) can be used as a marker of cardiovascular health status.
- Entropy measures represent a family of new methods to quantify the variability of the heart rate.
- Promising approach, due to ability to discover certain patterns and shifts in the "apparent ensemble amount of randomness" of stochastic processes,
- measure randomness and **predictability of processes**.

Mayer, C., Bachler, M., Holzinger, A., Stein, P. K. & Wassertheurer, S. 2016. The Effect of Threshold Values and Weighting Factors on the Association between Entropy Measures and Mortality after Myocardial Infarction in the Cardiac Arrhythmia Suppression Trial (CAST). *Entropy*, 18, (4), 129, doi:10.3390/e18040129.

Holzinger Group, hci-kdd.org

78

MAKE Health Verona 02



Mayer, C., Bachler, M., Holzinger, A., Stein, P. K. & Wassertheurer, S. 2016. The Effect of Threshold Values and Weighting Factors on the Association between Entropy Measures and Mortality after Myocardial Infarction in the Cardiac Arrhythmia Suppression Trial (CAST). *Entropy*, 18, (4), 129, doi:10.3390/e18040129.

Holzinger Group, hci-kdd.org

79

MAKE Health Verona 02

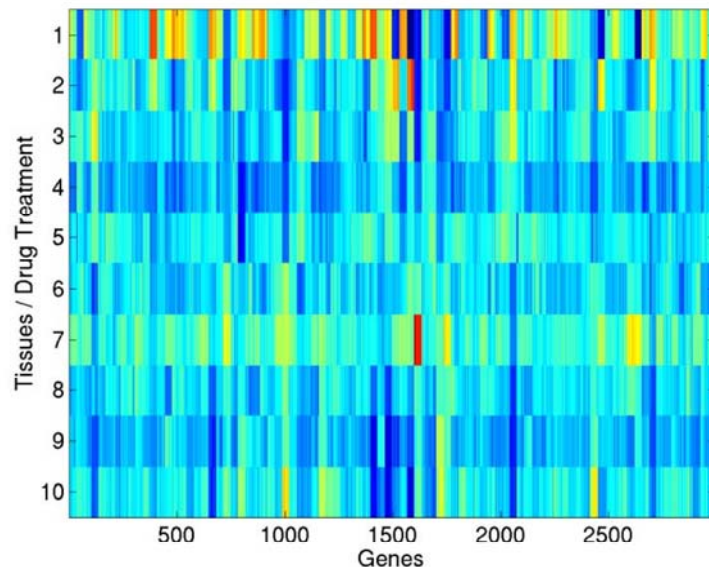
06 Cross-Entropy Kullback-Leibler Divergence

Holzinger Group, hci-kdd.org

80

MAKE Health Verona 02

- Entropy:
 - Measure for the **uncertainty** of random variables
- Kullback-Leibler divergence:
 - **comparing two distributions**
- Mutual Information:
 - measuring the **correlation** of two random variables



ON INFORMATION AND SUFFICIENCY

By S. KULLBACK AND R. A. LEIBLER

The George Washington University and Washington, D. C.

1. Introduction. This note generalizes to the abstract case Shannon's definition of information [15], [16]. Wiener's information (p. 75 of [18]) is essentially the same as Shannon's although their motivation was different (cf. footnote 1, p. 95 of [16]) and Shannon apparently has investigated the concept more completely. R. A. Fisher's definition of information (intrinsic accuracy) is well known (p. 709 of [6]). However, his concept is quite different from that of Shannon and Wiener, and hence ours, although the two are not unrelated as is shown in paragraph 2.

R. A. Fisher, in his original introduction of the *criterion of sufficiency*, required "that the statistic chosen should summarize the whole of the relevant information supplied by the sample," (p. 316 of [5]). Halmos and Savage in a recent paper, one of the main results of which is a generalization of the well known Fisher-Neyman theorem on sufficient statistics to the abstract case, conclude, "We think that confusion has from time to time been thrown on the subject by . . . , and (c) the assumption that a sufficient statistic contains all the information in only the technical sense of 'information' as measured by variance," (p. 241 of [8]). It is shown in this note that the information in a sample as defined herein, that is, in the Shannon-Wiener sense cannot be increased by any statistical operations and is invariant (not decreased) if and only if sufficient statistics are employed. For a similar property of Fisher's information see p. 717 of [6], Doob [19].

We are also concerned with the statistical problem of discrimination ([3], [17]), by considering a measure of the "distance" or "divergence" between statistical populations ([1], [2], [13]) in terms of our measure of information. For the statistician two populations differ more or less according as to how difficult it is to discriminate between them with the best test [14]. The particular measure of divergence we use has been considered by Jeffreys ([10], [11]) in another connection. He is primarily concerned with its use in providing an invariant density of *a priori* probability. A special case of this divergence is Mahalanobis' generalized distance [13].



Solomon Kullback
1907-1994



Richard Leibler
1914-2003

Kullback, S. & Leibler, R. A. 1951. On information and sufficiency. The annals of mathematical statistics, 22, (1), 79-86, www.jstor.org/stable/2236703

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

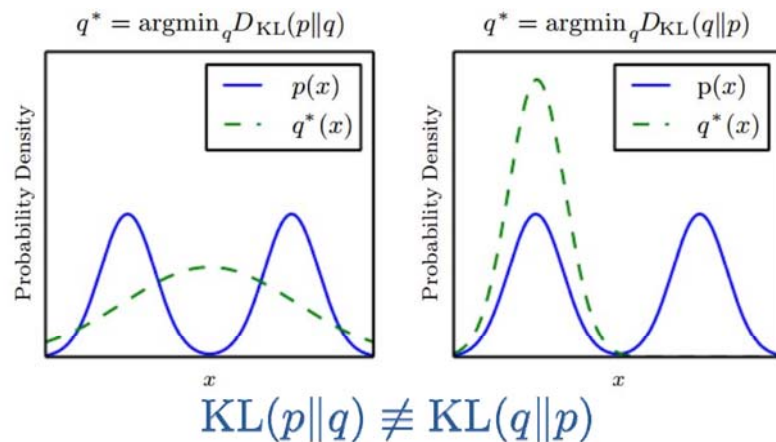
Shannon, C. E. 1948. A Mathematical Theory of Communication. Bell System Technical Journal, 27, 379-423.

Important quantity in

- coding theory
- statistical physics
- machine learning

$$H[y|x] = - \iint p(y, x) \ln p(y|x) dy dx$$

$$H[x, y] = H[y|x] + H[x]$$



Goodfellow, I., Bengio, Y. & Courville, A. 2016. Deep Learning, Cambridge (MA), MIT Press.

$$KL(p||q) = - \int p(x) \ln q(x) dx - \left(- \int p(x) \ln p(x) dx \right)$$

$$= - \int p(x) \ln \left\{ \frac{q(x)}{p(x)} \right\} dx$$

$$KL(p||q) \simeq \frac{1}{N} \sum_{n=1}^N \{ - \ln q(x_n|\theta) + \ln p(x_n) \}$$

$$KL(p||q) \geq 0$$

KL-divergence is often used to measure the distance between two distributions

- ... are **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**

Mutual Information and Point Wise MI

Example: Disease-Disease Relationship

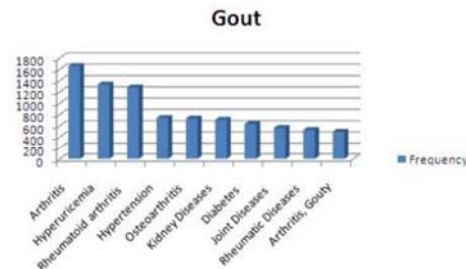
Let two words, w_i and w_j , have probabilities $P(w_i)$ and $P(w_j)$.
Then their mutual information $PMI(w_i, w_j)$ is defined as:

$$PMI(w_i, w_j) = \log \left(\frac{P(w_i, w_j)}{P(w_i) P(w_j)} \right)$$

For w_i denoting *rheumatoid arthritis* and w_j representing *diffuse scleritis* the following simple calculation yields:

$$P(w_i) = \frac{94,834}{20,033,079} \quad P(w_j) = \frac{74}{20,033,079}$$

$$P(w_i, w_j) = \frac{13}{94,834} \quad PMI(w_i, w_j) = 7.7.$$



$$I[\mathbf{x}, \mathbf{y}] \equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y}))$$

$$= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

- Measures how much reduction in uncertainty of X given the information about Y
- Measures correlation between X and Y
- Related to the “channel capacity” in the original Shannon information theory

Bishop, C. M. 2007. *Pattern Recognition and Machine Learning*, Heidelberg, Springer.

$$SCP(x, y) = p(x|y) \cdot p(y|x) = \frac{p(x, y)}{p(y)} \cdot \frac{p(x, y)}{p(x)} = \frac{p(x, y)^2}{p(x) \cdot p(y)}$$

Table 4 Comparison of FACTAs ranking of related concepts from the category Symptom for the query “rheumatoid arthritis” created by the methods co-occurrence frequency, PMI, and SCP

Frequency		PMI		SCP	
pain	5667	impaired body balance	7,8	swollen joints	0.002
Arthralgia	661	ASPIRIN INTOLERANCE	7,8	pain	0.001
fatigue	429	Epitrochlear lymphadenopathy	7,8	Arthralgia	0.001
diarrhea	301	swollen joints	7,4	fatigue	0.000
swollen joints	299	Joint tenderness	7	erythema	0.000
erythema	255	Occipital headache	6,2	splenomegaly	0.000
Back Pain	254	Neuromuscular excitation	6,2	Back Pain	0.000
headache	239	Restless sleep	5,8	polymyalgia	0.000
splenomegaly	228	joint crepitus	5,7	joint stiffness	0.000
Anesthesia	221	joint symptom	5,5	Joint tenderness	0.000
dyspnea	218	Painful feet	5,5	hip pain	0.000
weakness	210	feeling of malaise	5,5	metatarsalgia	0.000
nausea	199	Homan's sign	5,4	Skin Manifestations	0.000
Recovery of Function	193	Diffuse pain	5,2	neck pain	0.000
low back pain	167	Palmar erythema	5,2	Eye Manifestations	0.000
abdominal pain	141	Abnormal sensation	5,2	low back pain	0.000

- 1) Challenges include –omics data analysis, where KL divergence and related concepts could provide important **measures** for discovering biomarkers.
- 2) Hot topics are new entropy measures suitable for computations in the context of complex/uncertain data for ML algorithms.
- Inspiring is the abstract geometrical setting underlying ML main problems, e.g. Kernel functions can be completely understood in this perspective. Future work may include entropic concepts and geometrical settings.



Thank you!

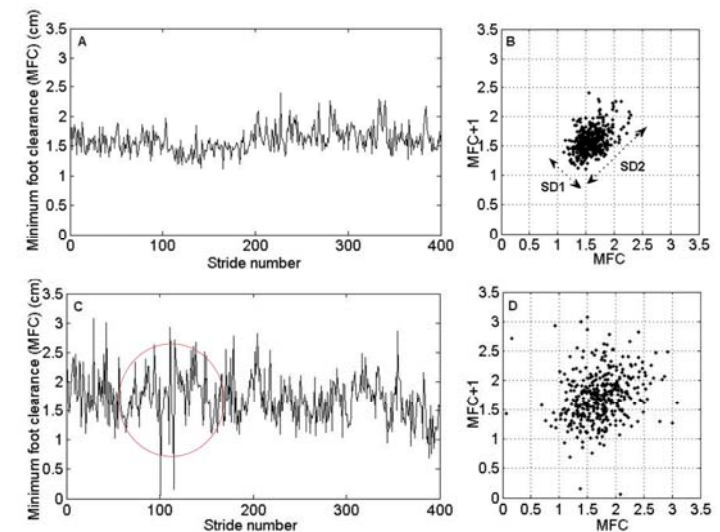
- The case of higher order statistical structure in the data – nonlinear and hierarchical ?
- Outliers in the data – noise models?
- There are $\frac{D(D+1)}{2}$ parameters in a multi-variate Gaussian model – what happens if $D \gg ?$
dimensionality reduction

Questions

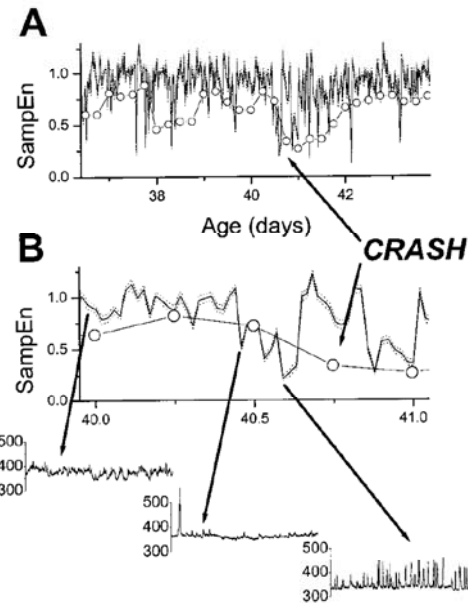
- What are the grand challenges in ML for health?
- What is the key problem before you can apply ML?
- Describe the taxonomy of data at Hospital level!
- What does translational medicine mean?
- Give an example for a 2.5D-data set!
- Why would be the combination of ontologies with machine learning provide a benefit?
- How did Van Bommel and Musen describe the interplay between data-information-knowledge?
- What is the “body-of-knowledge” in medical jargon?
- How do human process information?

Appendix

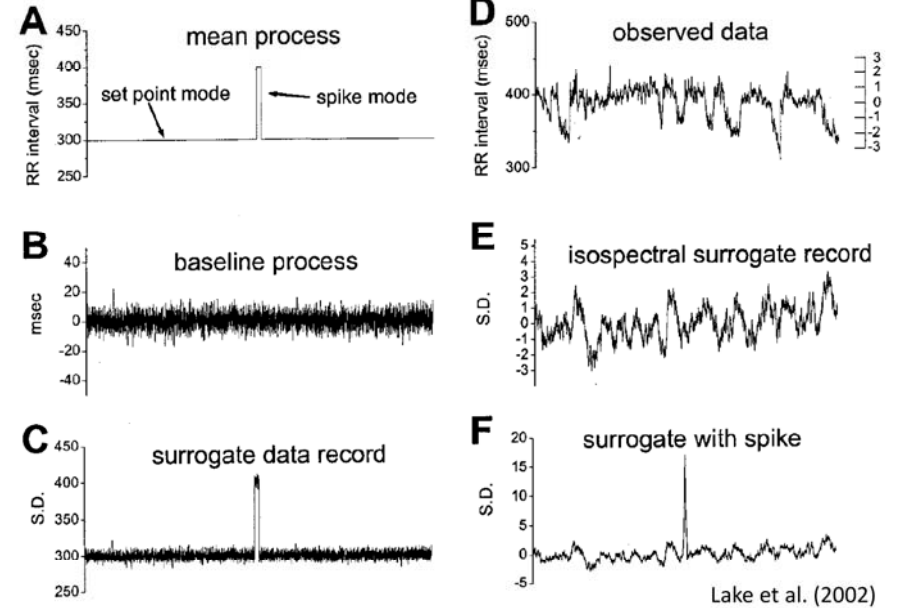
- What was our definition of “knowledge”?
- What is the huge benefit of a probabilistic model?
- Please explain Bayes law with view on ML!
- What is information in the sense of Shannon?
- Why is information theory for us important?
- Which benefits provide entropic methods for us?
- Why is feature selection so important?
- What can you do with the Kullback-Leibler Divergence?



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.



Lake, D. E., Richman, J. S., Griffin, M. P. & Moorman, J. R. (2002) Sample entropy analysis of neonatal heart rate variability. *American Journal of Physiology-Regulatory Integrative and Comparative Physiology*, 283, 3, R789-R797.



Lake et al. (2002)

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 [1]:

- 1) Form m -vectors, $X(i)$ to $X(N-m+1)$ defined by:

$$X(i) = [x(i), x(i+1), \dots, x(i+m-1)] \quad i = 1, N-m+1 \quad (1)$$
- 2) Define the distance $d[X(i), X(j)]$ between vectors $X(i)$ and $X(j)$ as the maximum absolute difference between their respective scalar components:

$$d[X(i), X(j)] = \max_{k=0, m-1} [|x(i+k) - x(j+k)|] \quad (2)$$
- 3) Define for each i , for $i=1, N-m+1$, let

$$C_r^m(i) = V^m(i) / (N-m+1) \quad (3)$$

where $V^m(i) = \text{no. of } d[X(i), X(j)] \leq r$
- 4) Take the natural logarithm of each $C_r^m(i)$, and average it over i as defined in step 3):

$$\phi^m(r) = \frac{1}{N-m+1} \sum_{i=1}^{N-m+1} \ln(C_r^m(i)) \quad (4)$$
- 5) Increase the dimension to $m+1$ and repeat steps 1) to 4).
- 6) Calculate ApEn value for a finite data length of N :

$$\text{ApEn}(m, r, N) = \phi^m(r) - \phi^{m+1}(r) \quad (5)$$

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

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, $X(i)$ to $X(N-m+1)$ defined by:

$$X(i) = [x(i), x(i+1), \dots, x(i+m-1)] \quad i = 1, N-m+1 \quad (6)$$
- 2) Define the distance $d_m[X(i), X(j)]$ between vectors $X(i)$ and $X(j)$ as the maximum absolute difference between their respective scalar components:

$$d_m[X(i), X(j)] = \max_{k=0, m-1} [|x(i+k) - x(j+k)|] \quad (7)$$
- 3) Define for each i , for $i=1, N-m$, let

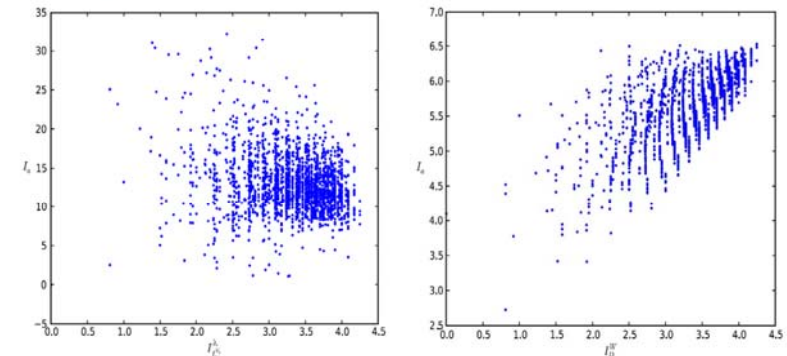
$$B_i^m(r) = \frac{1}{N-m-1} \times \text{no. of } d_m[X(i), X(j)] \leq r, i \neq j \quad (8)$$
- 4) Similarly, define for each j , for $j=1, N-m$, let

$$A_j^m(r) = \frac{1}{N-m-1} \times \text{no. of } d_{m+1}[X(i), X(j)] \leq r, i \neq j \quad (9)$$
- 5) Define $B^m(r) = \frac{1}{N-m} \sum_{i=1}^{N-m} B_i^m(r)$ (10)

$$A^m(r) = \frac{1}{N-m} \sum_{j=1}^{N-m} A_j^m(r) \quad (11)$$
- 6) SampEn value for a finite data length of N can be estimated:

$$\text{SampEn}(m, r, N) = -\ln(A^m(r)/B^m(r)) \quad (12)$$

- 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.

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>

- **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);
- **Artefact/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;

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>

- **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

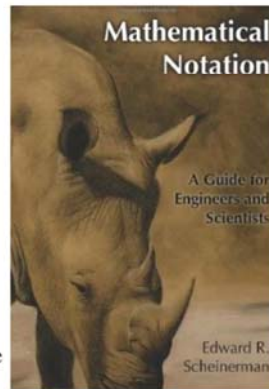
*"In mathematics you don't understand things. You just get used to them" –
John von Neumann*

Data

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

Distribution

P	Probability
$F(\mathbf{x})$	Cumulative probability distribution function (cdf)
$p(\mathbf{x})$	Probability density function (pdf)
$p(\mathbf{x}, \mathbf{y})$	Joint probability density function
$p(\mathbf{x}; \omega)$	Probability density function, which is parameterized
$p(\mathbf{y} \mathbf{x})$	Conditional density
$t(\mathbf{x})$	Target function

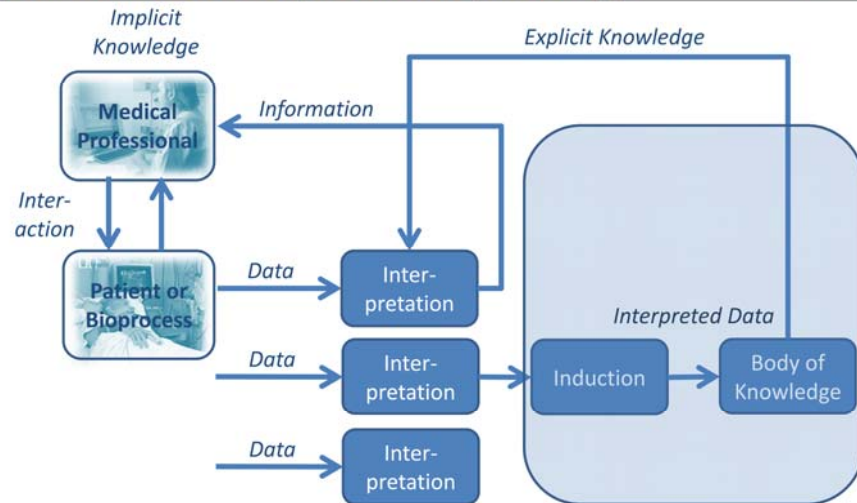


Scientists in data integration - selection - incomplete!

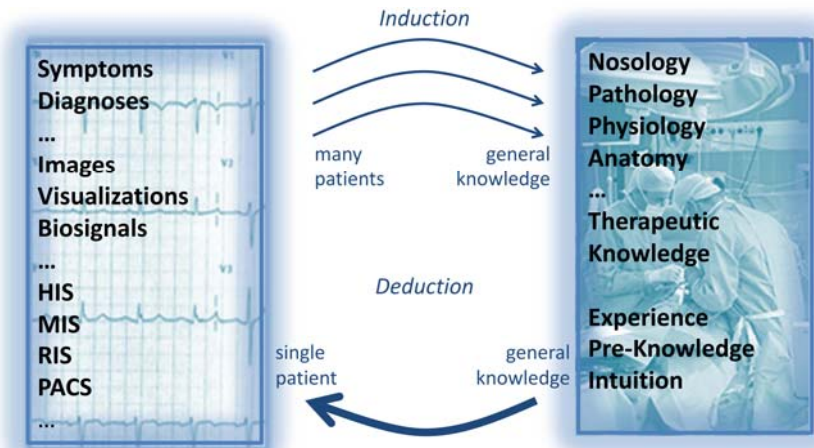
Status as of 04.04.2016

- ApEn = Approximate Entropy;
- \mathbb{C}_{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}_{data} = Data in perceptual space;
- PDB = Protein Data Base;
- SampEn = Sample Entropy;

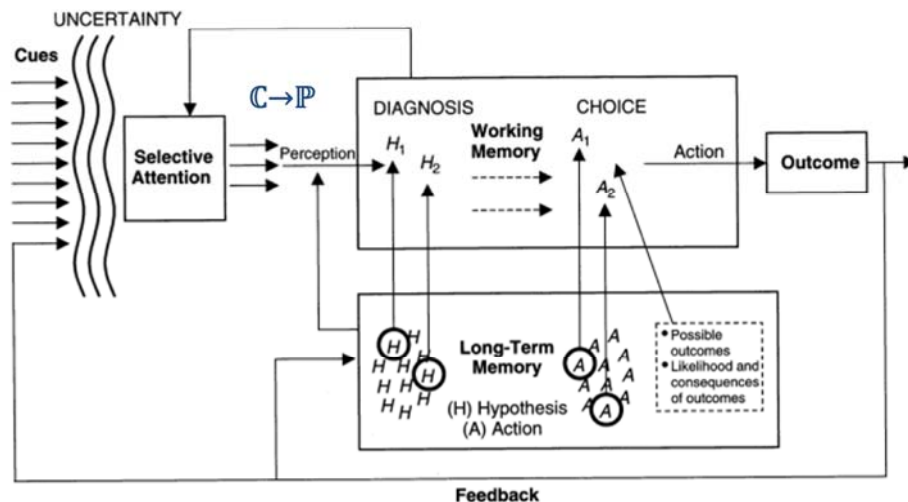
Clinical view on data – information, and knowledge



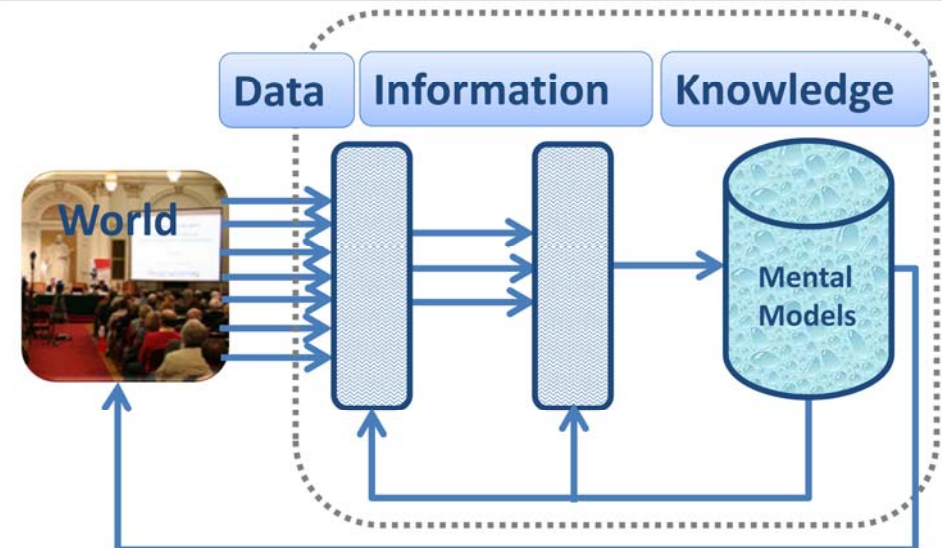
Bemmel, J. H. v. & Musen, M. A. (1997) *Handbook of Medical Informatics*. Heidelberg, Springer.



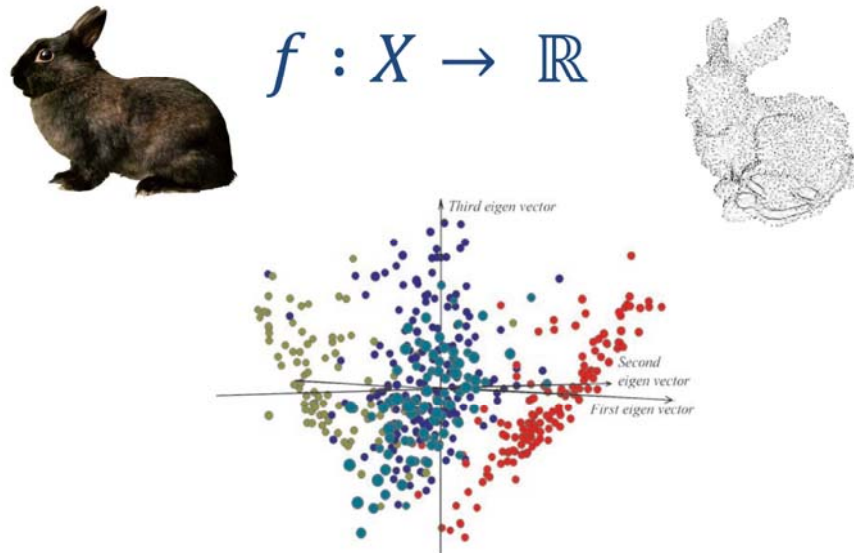
Holzinger (2007)



Wickens, C. D. (1984) *Engineering psychology and human performance*. Columbus: Merrill.



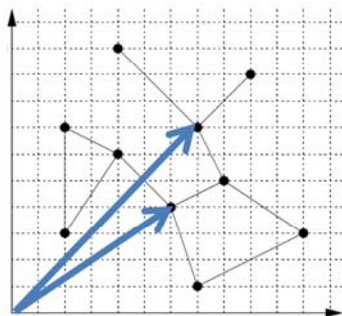
Knowledge := a set of expectations



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.

Example Metric Space

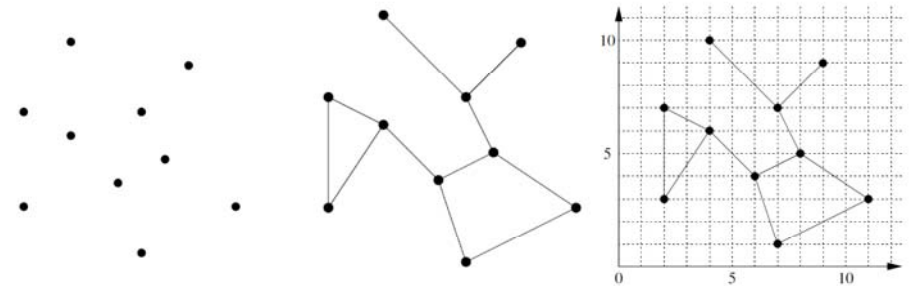
A set S with a metric function d is a metric space



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

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

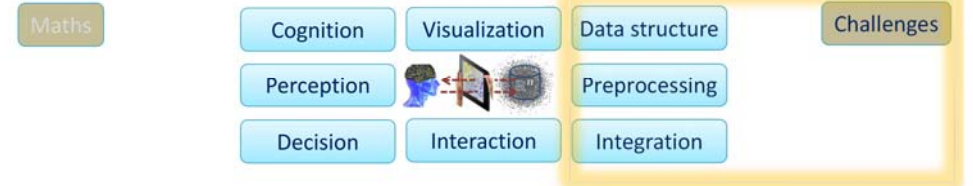
Let us collect n -dimensional i observations: $x_i = [x_{i1}, \dots, x_{in}]$



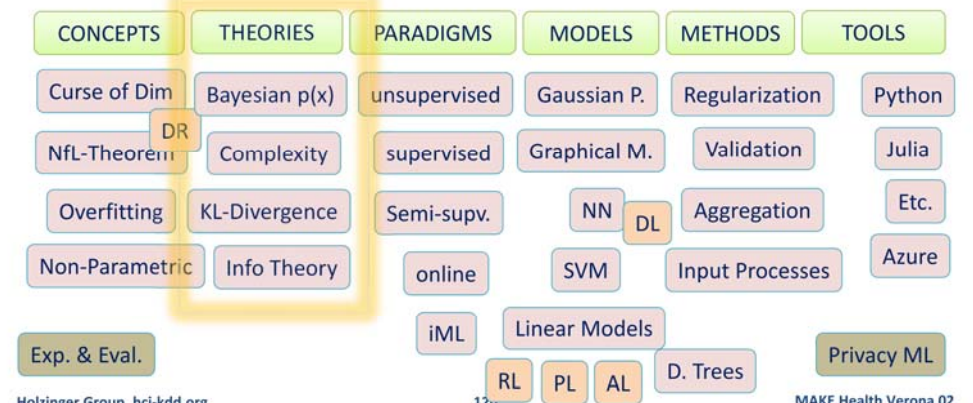
Point cloud in \mathbb{R}^2 topological space metric space

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

ML-Jungle Top Level View and the focus of today ...



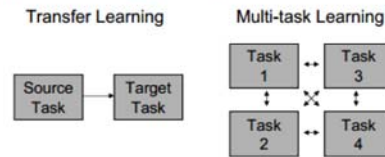
Always with a focus/application in health informatics



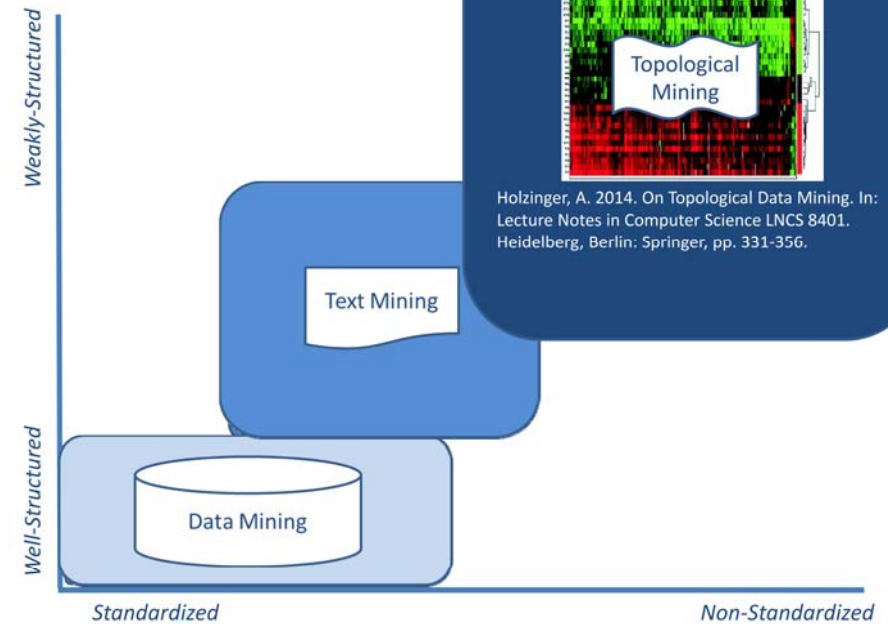
- Big data with many training sets (this is good for ML!)
- Small number of data sets, rare events
- Very-high-dimensional problems
- Complex data – NP-hard problems
- Missing, dirty, wrong, noisy, ..., data

GENERALISATION

TRANSFER



Torrey, L. & Shavlik, J. 2009. Transfer learning. Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques, 242-264, doi:10.4018/978-1-60566-766-9.ch011.



Holzinger, A. 2014. On Topological Data Mining. In: Lecture Notes in Computer Science LNCS 8401. Heidelberg, Berlin: Springer, pp. 331-356.

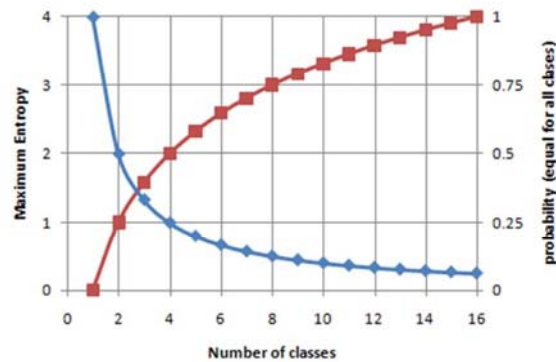
- $X: S \rightarrow \mathbb{R}$ ("measure" of outcome)
- Events can be defined according to X
 - $E(X=a) = \{s_i | X(s_i)=a\}$
 - $E(X \geq a) = \{s_i | X(s_i) \geq a\}$
- Consequently, probabilities can be defined on X
 - $P(X=a) = P(E(X=a))$
 - $P(a \geq X) = P(E(a \geq X))$
- partitioning the sample space



My greatest concern was what to call it. I thought of calling it "information", but the word was overly used, so I decided to call it "uncertainty". When I discussed it with John von Neumann, he had a better idea. Von Neumann told me, "You should call it entropy, for two reasons. In the first place your uncertainty function has been used in statistical mechanics under that name, so it already has a name. In the second place, and more important, nobody knows what entropy really is, so in a debate you will always have the advantage."

Tribus, M. & McIrvine, E. C. (1971) Energy and Information. *Scientific American*, 225, 3, 179-184.

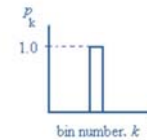
$$\log_2 \frac{1}{p} = -\log_2 p \quad H = -\sum_{i=1}^N p_i \log_2(p_i)$$



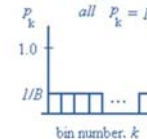
Shannon, C. E. (1948) A Mathematical Theory of Communication. *Bell System Technical Journal*, 27, 379-423.

- Developed by Claude Shannon in the 1940s
- Maximizing the amount of information that can be transmitted over an imperfect communication channel
- Data compression (entropy)
- Transmission rate (channel capacity)

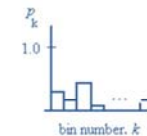
Claude E. Shannon: *A Mathematical Theory of Communication*, *Bell System Technical Journal*, Vol. 27, pp. 379-423, 623-656, 1948



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



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



$$H = H_{max} = \log_2 N$$



vapnik

Professor of Columbia, Fellow of NEC Labs America, machine learning, statistics, computer science
Verified email at nec-labs.com

Title 1-20

The Nature of Statistical Learning Theory

V Vapnik
Data mining and knowledge discovery

Follow

Google Scholar

Get my own profile

Citation indices	All	Since 2012
Citations	184842	78672
h-index	116	76
i10-index	386	277

- The VC dimension is a measure of the capacity of a space of functions that can be learned by a statistical classification algorithm. It is defined as the cardinality of the largest set of points that the algorithm can shatter. It is a core concept in Vapnik-Chervonenkis theory

Vapnik, V. N. & Chervonenkis, A. Y. 1971. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability & Its Applications*, 16, (2), 264-280, doi:10.1137/1116025.

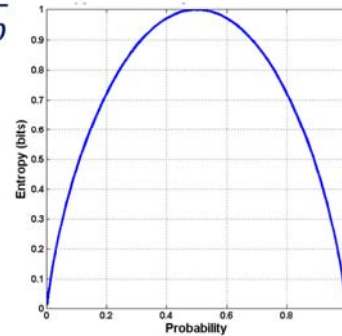
$$Q \dots P = \{p_1, \dots, p_n\} \quad H(Q) = - \sum_{i=1}^n (p_i * \log p_i)$$

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

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

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.



- 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, for modeling etc. 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.