


Andreas Holzinger
185.A83 Machine Learning for Health Informatics
2017S, VU, 2.0 h, 3.0 ECTS
Lecture 06 - Module 04 – Week 17 - 25.04.2017

Probabilistic Graphical Models
Part 2: From Bayesian Networks
Probabilistic Topic Models


a.holzinger@hci-kdd.org
<http://hci-kdd.org/machine-learning-for-health-informatics-course>



Holzinger Group, hci-kdd.org 1 Machine Learning Health 06

Red thread through the lecture today

- 01 Probabilistic Decision Making
- 02 Probabilistic Topic Models
- 03 Knowledge Representation in Net Medicine
- 04 ML on Graphs Examples
- 05 Digression: Similarity
- 06 Graph Measures
- 07 Point Clouds from Natural Images



Holzinger Group, hci-kdd.org 4 Machine Learning Health 06

A fundamental problem first ...

$$\mathbb{E}[f] = \int f(z)p(z)dz$$

↓

$$\hat{f} = \frac{1}{L} \sum_{l=1}^L f(z^{(l)})$$

Holzinger Group, hci-kdd.org 7 Machine Learning Health 06

ML needs a concerted effort fostering integrated research

<http://hci-kdd.org/international-expert-network>

Interactive Data Mining Knowledge Discovery

6 Data Visualization
2 Learning Algorithms
1 Data Mapping Pre-processing Data Fusion

3 GDM Graph-based Data Mining
4 TDM Topological Data Mining
5 EDM Entropy-based Data Mining

7 Privacy, Data Protection, Safety and Security

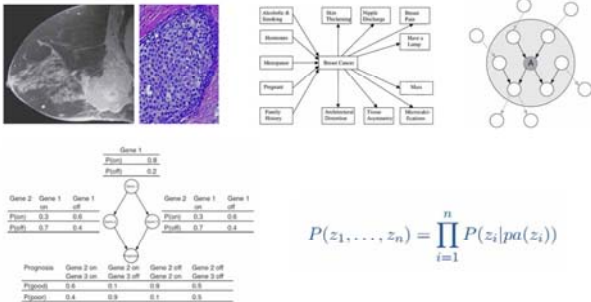
Holzinger, A. 2014. Trends in Interactive Knowledge Discovery for Personalized Medicine: Cognitive Science meets Machine Learning. IEEE Intelligent Informatics Bulletin, 15, (1), 6-14.

Holzinger Group, hci-kdd.org 2 Machine Learning Health 06



Holzinger Group, hci-kdd.org 5 Machine Learning Health 06

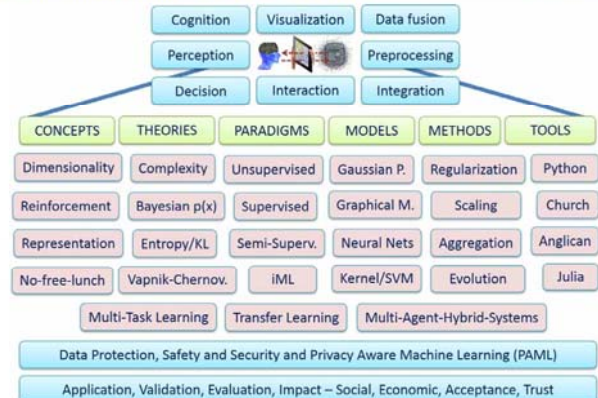
Medical Example: Breast cancer prognosis incl. Genetics



Gevaert, O., Smet, F. D., Timmerman, D., Moreau, Y. & Moor, B. D. (2006) Predicting the prognosis of breast cancer by integrating clinical and microarray data with Bayesian networks. Bioinformatics, 22, 14, 184-190.

Holzinger Group, hci-kdd.org 8 Machine Learning Health 06

Machine Learning Jungle Top-Level View



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 3 Machine Learning Health 06

To reach a level of usable intelligence we need to ...

- 1) learn from prior data
- 2) extract knowledge
- 2) generalize,
 - i.e. guessing where a probability mass function concentrates
- 4) fight the curse of dimensionality
- 5) disentangle underlying explanatory factors of data, i.e.
- 6) understand the data in the context of an application domain

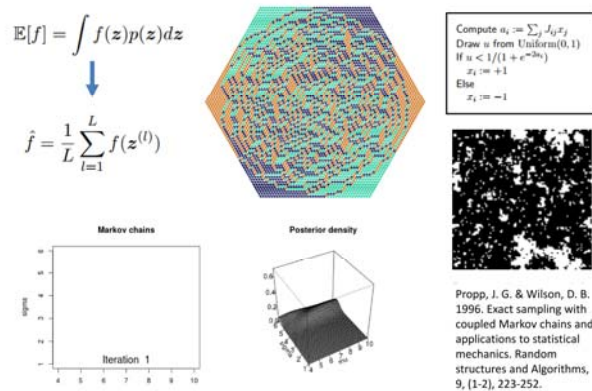
Holzinger Group, hci-kdd.org 6 Machine Learning Health 06

Inference in Bayes Nets is intractable (NP-complete!)

- For certain cases it is tractable if:
 - Just one variable is unobserved
 - We have singly connected graphs (no undirected loops -> belief propagation)
 - Assigning probability to fully observed set of variables
- Possibility: Monte Carlo Methods (generate many samples according to the Bayes Net distribution and then count the results)
- Otherwise: approximate solutions, NOTE:

Sometimes it is better to have an approximate solution to a complex problem – than a perfect solution to a simplified problem

Holzinger Group, hci-kdd.org 9 Machine Learning Health 06



"I saw her duck"



Radiologischer Befund

angelegt am 06.05.2006 09:26
geprüft von ...
gedruckt am 17.11.2006 08:24
Abk.: NCHB

Kurzname: St.p. SHT
Fragestellung: -
Untersuchung: Thorax eine Ebene liegend

SB
Bewegungsartefakte. Zustand nach Schädelhirntrauma.
Das Cor in der Größenform, keine akuten Stauungszeichen.
Fragliches Infiltrat paravertebral li. im UF, RW-Erguss li.
Zustand nach Anlage eines ET, die Spitze ca. 5cm cranial der Bifurkation, liegt MS, orthotop positioniert. ZVK über re., die Spitze in Proj. auf die VCS. Kein Hinweis auf Pneumothorax.
Der re. Rezessus frei.

Mit kollegialen Grüßen

Elektronische Freigabe durch ... am 06.05.2006

**Special Words
Language Mix
Abbreviations
Errors ...**

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

- HWI =
 - Harnwegsinfekt
 - Hinterwandinfarkt
 - Hinterwandischämie
 - Hakenwurminfektion
 - Halswirbelimmobilisation
 - Hip Waist Index
 - Height-Width Index
 - Heart-Work Index
 - Hemodynamically weighted imaging
 - High Water Intake
 - Hot water irrigation
 - Hepatic weight index
 - Häufig wechselnder Intimpartner



- Leitung = Nervenleitung, Abteilungsleitung, Stromleitung, Wasserleitung, Harnleitung, Ableitung, Vereinsleitung ☺...

01 Probabilistic Decision Making

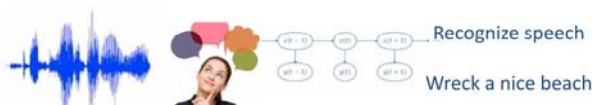
Laplace, P.-S. 1781. Mémoire sur les probabilités. *Mémoires de l'Académie Royale des sciences de Paris*, 1778, 227-332.



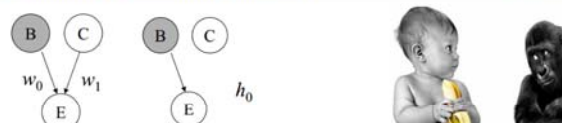
- Example 1: Inverse Probability
- Example 2: Diagnosis
- Example 3: Language understanding:

$$p(h|d) \propto p(D|\theta) * p(h)$$

$$P(words|sounds) \propto P(sounds|words) * P(words)$$



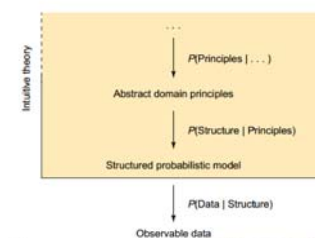
- Learning ensures that new observations (d) match our previous hypotheses (h)

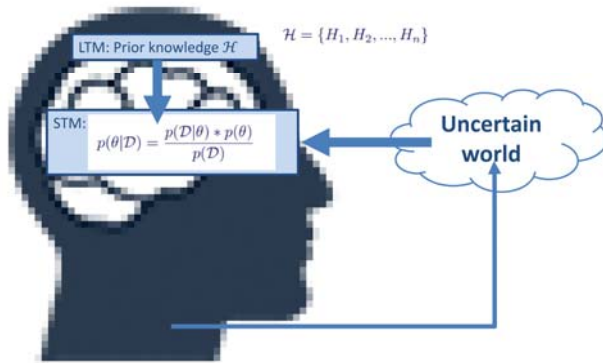


- Visual perception, language understanding, motor learning, associative learning, categorization, concept learning, reasoning, causal inference, ...
- Learning concepts from (few!) examples
- Learning and applying intuitive theories (balancing complexity vs. fit optimality)

- Similarity
- Representativeness and evidential support
- Causal judgement
- Coincidences and causal discovery
- Diagnostic inference
- Predicting the future

Tenenbaum, J. B., Griffiths, T. L. & Kemp, C. 2006. Theory-based Bayesian models of inductive learning and reasoning. *Trends in cognitive sciences*, 10, (7), 309-318.





Expected Utility Theory $E(U|d)$

For a single decision variable an agent can select $D = d$ for any $d \in \text{dom}(D)$.

The expected utility of decision $D = d$ is



<http://www.eohh.info/page/Oskar+Morgenstern>

$$E(U | d) = \sum_{x_1, \dots, x_n} P(x_1, \dots, x_n | d) U(x_1, \dots, x_n, d)$$

An optimal single decision is the decision $D = d_{\max}$ whose expected utility is maximal:

$$d_{\max} = \arg \max_{d \in \text{dom}(D)} E(U | d)$$

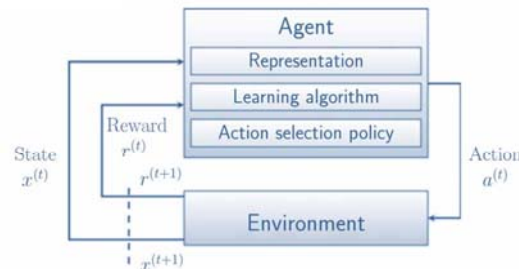
Von Neumann, J. & Morgenstern, O. 1947. Theory of games and economic behavior, Princeton university press.

Case	Case ID	Case Name	Case Type	Case Status	Case Date	Case Location	Case Description	Case Notes	Case Actions	Case Results	Case Comments
2010018065	ambulatorischer Fall	13.01.2010	MO-KARDIO MK-KardioAmb								
	Leistungen (VWL, RAD, Therapie)										
	EXG (12 Ableitungen)	13.01.2010	08:00	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2010018065	
	HR-Intervall-Untersuchung	13.01.2010	08:00	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2010018065	
	Schrittmacherkontrolle	13.01.2010	08:00	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2010018065	
2010022197	ambulatorischer Fall	04.01.2010	CH-TRANSF CK-Transp								
	Diagnosen Gesamt (3)										
2009494995	stationärer Fall	20.12.2009	MEDANNO Med Angio								
	Diagnosen Gesamt (14)										
	Leistungen (VWL, RAD, Therapie)										
	Becken-u. Beckenboden	22.12.2009	16:38	RIVIRACB RVH-Raum B	STANNELI	OK				2009494995	
	Laufbandergometer	21.12.2009	08:30	MA-ANGIO MK-AngioAmb	SPARANDR	OK				2009494995	
	Herzrhythmusuntersuchung	21.12.2009	08:30	MA-ANGIO MK-AngioAmb	SPARANDR	OK				2009494995	
2009453621	stationärer Fall	17.11.2009	CK-TX-MC								
	Diagnosen Gesamt (12)										
	Leistungen (VWL, RAD, Therapie)										
	Therapie (12 Ableitungen)	23.11.2009	08:05	CK-PHYVIO CK-Phevio	BEITWALT	OK				2009453621	
	Organop. Bildwandergel	17.11.2009	08:12	CK-TX-OP	SCHWABCH	OK				2009453621	
2009431136	ambulatorischer Fall	28.10.2009	MO-KARDIO MK-KardioAmb								
	Leistungen (VWL, RAD, Therapie)										
	Schrittmacherkontrolle	28.10.2009	09:15	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2009431136	
	HR-Intervall-Untersuchung	28.10.2009	09:15	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2009431136	
	EXG (12 Ableitungen)	28.10.2009	09:15	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2009431136	
	Fluoridokumentation, Video	28.10.2009	09:15	MO-KARDIO MK-KardioAmb	DUELTIMO	OK				2009431136	
2009378733	ambulatorischer Fall	18.09.2009	MO-KARDIO MK-KardioAmb								
	Diagnosen Gesamt (5)										
	Leistungen (VWL, RAD, Therapie)										
	Rudbeck-Largeit	24.09.2009	10:59	MO-KARDIO MK-KardioAmb	RUDBECK	OK				2009378733	
	Rudbeck-Largeit	24.09.2009	12:02	MO-KARDIO MK-KardioAmb	RUDBECK	OK				2009378733	
2009187548	stationärer Fall	21.04.2009	CK-OM-K								
	Diagnosen Gesamt (5)										
	Leistungen (VWL, RAD, Therapie)										
	Fotodokumentation, Video	28.04.2009	08:49	MO-KARDIO MK-KardioAmb	PITTHEID	OK				2009187548	
	EXG (12 Ableitungen)	28.04.2009	08:49	MO-KARDIO MK-KardioAmb	PITTHEID	OK				2009187548	
	HR-Intervall-Untersuchung	28.04.2009	08:49	MO-KARDIO MK-KardioAmb	PITTHEID	OK				2009187548	
	Schrittmacherkontrolle	28.04.2009	08:49	MO-KARDIO MK-KardioAmb	PITTHEID	OK				2009187548	
	EXG (12 Ableitungen)	23.04.2009	10:43	MO-KARDIO MK-KardioAmb	KOBENROR	OK				2009187548	
	HR-Intervall-Untersuchung	23.04.2009	10:43	MO-KARDIO MK-KardioAmb	KOBENROR	OK				2009187548	
	Kardiografie	21.04.2009	10:22	MO-KARDIO MK-KardioAmb	LANAMICH	OK				2009187548	

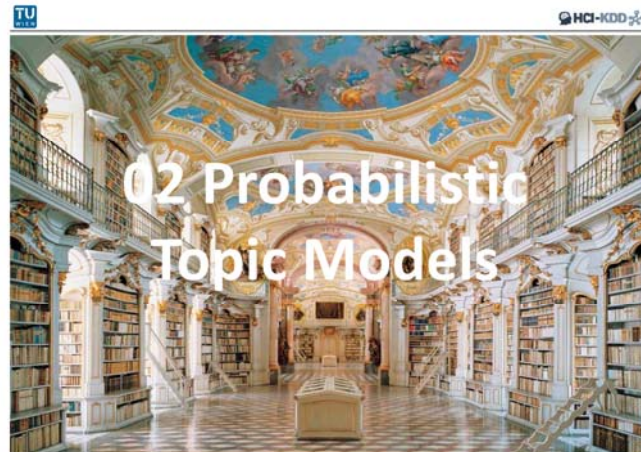
RL-Agent seeks to maximize rewards

```
for t = 1, ..., n do
  The agent perceives state s.
  The agent performs action a.
  The environment evolves to s'.
  The agent receives reward r.
end for
```

Intelligent behavior arises from the actions of an individual seeking to maximize its received reward signals in a complex and changing world



Sutton, R. S. & Barto, A. G. 1998. Reinforcement learning: An introduction, Cambridge MIT press



Example (1)

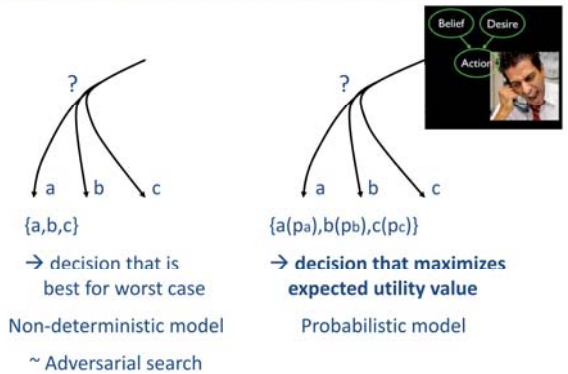
$$D = \langle d_1, d_2, \dots, d_n \rangle$$

$$d_i = t_1, t_2, \dots, t_k$$

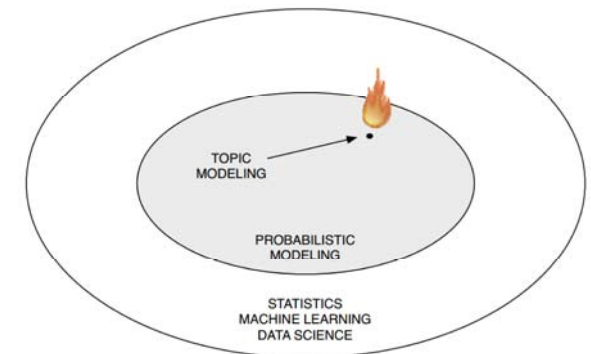
$$w_{i,j} = \begin{cases} 1, & t_i \in d_j \\ 0, & t_i \notin d_j \end{cases} \rightarrow d_j = (0, 1, 1, 0, 1, \dots, 1)^T$$

$$w_{i,j} = \begin{cases} (1 + \log f_{i,j}) * \log \frac{N}{n_i}, & \text{if } f_{i,j} > 0 \\ 0, & \text{otherwise} \end{cases}$$

De-cision (Ent-scheidung) between alternatives

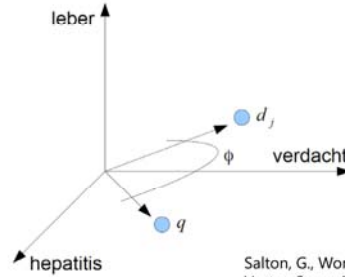
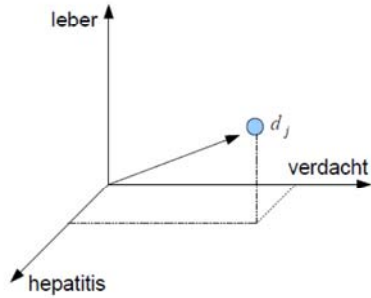


Topic modelling – small topic but hot topic in ML



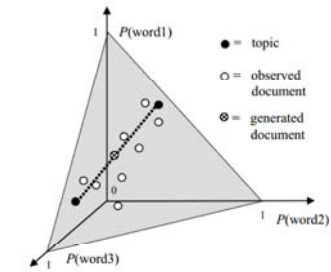
Example (2)

$$D_{m \times n} = \begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,n-1} & w_{1,n} \\ w_{2,1} & w_{2,2} & & w_{2,n-1} & w_{2,n} \\ \vdots & & & & \vdots \\ w_{m-1,1} & w_{m-1,2} & & w_{m-1,n-1} & w_{m-1,n} \\ w_{m,1} & w_{m,2} & \dots & w_{m,n-1} & w_{m,n} \end{bmatrix}$$



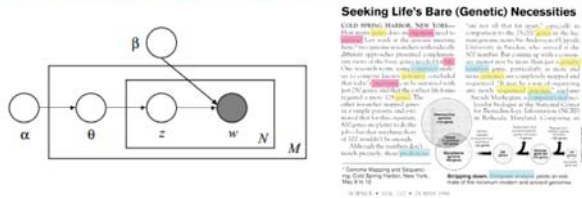
$$\cos(\phi) = \frac{q \cdot d_j}{\|q\| \|d_j\|}$$

Salton, G., Wong, A. & Yang, C. S. 1975. Vector-Space Model for automatic indexing. *Communications of the ACM*, 18, (11), 613-620.



- Documents = categorical distributions over a large space of predefined vocabulary
- Topics = categorical distributions
- Generative model = each document can be seen as a convex combination of the topic distributions

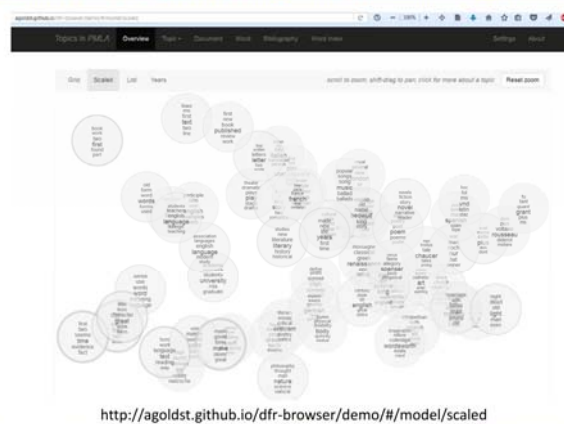
Teh, Y. W., Jordan, M. I., Beal, M. J. & Blei, D. M. 2006. Hierarchical dirichlet processes. *Journal of the american statistical association*, 101, (476), 1566-1581.



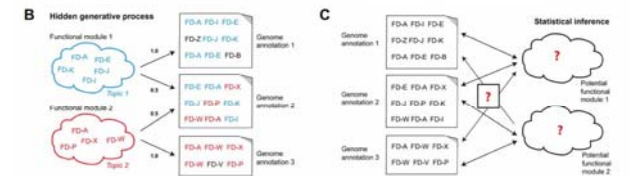
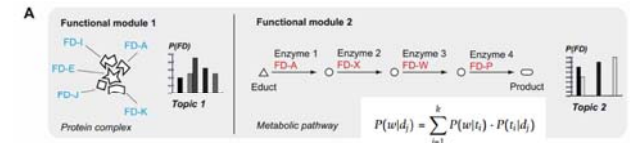
Given the parameters α and β , the joint distribution of a topic mixture θ , a set of N topics z , and a set of N words w is given by:

$$p(\theta, z, w | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

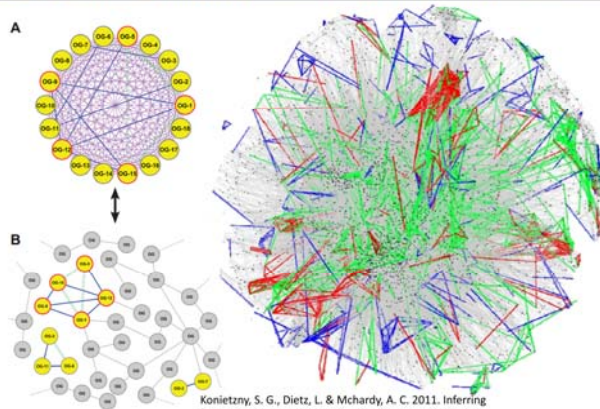
Blei, D. M., Ng, A. Y. & Jordan, M. I. 2003. Latent dirichlet allocation. *The Journal of machine Learning research*, 3, 993-1022.



<http://agoldst.github.io/dfr-browser/demo/#/model/scaled>

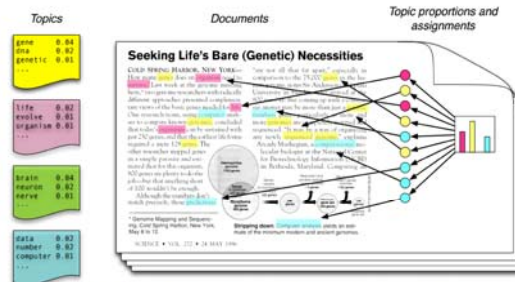


Konietzny, S. G., Dietz, L. & Mchardy, A. C. 2011. Inferring functional modules of protein families with probabilistic topic models. *BMC bioinformatics*, 12, (1), 1.

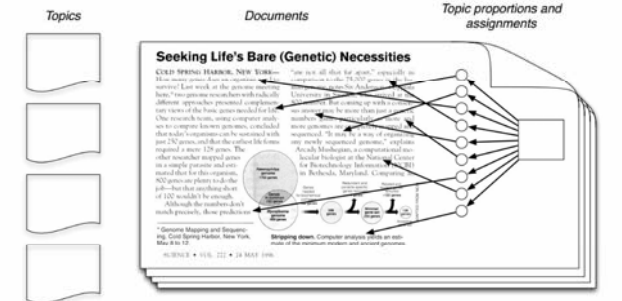


Konietzny, S. G., Dietz, L. & Mchardy, A. C. 2011. Inferring functional modules of protein families with probabilistic topic models. *BMC bioinformatics*, 12, (1), 1.

Goal: to get insight in unknown document collections
See a nice demo <http://agoldst.github.io/dfr-browser/demo/#/model/grid>

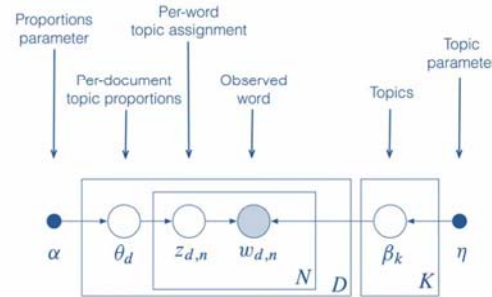


Each doc is a random mix of corpus-wide topics and each word is drawn from one of these topics

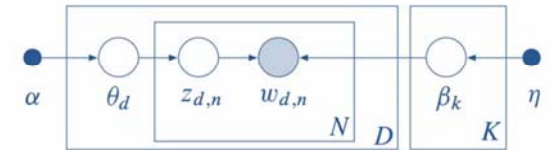


We only observe the docs – the other structure is hidden; then we compute the posterior $p(t,p,a|docs)$

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

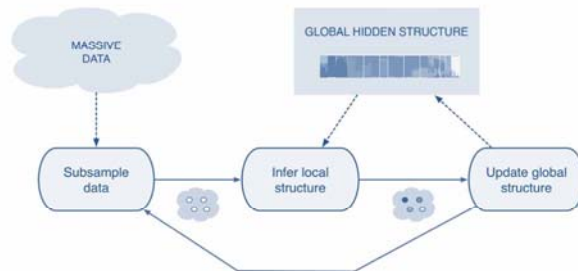


- Encodes assumptions on data with a factorization of the joint
- Connects assumptions to algorithms for computing with data
- Defines the posterior (through the joint)

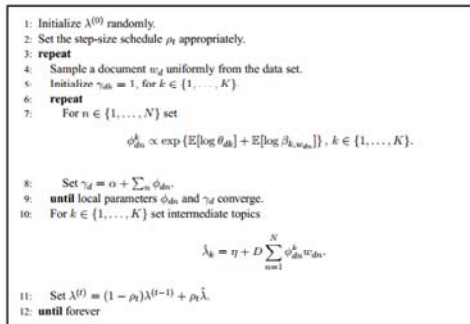


$$p(\beta, \theta, z | w) = \frac{p(\beta, \theta, z, w)}{\int_{\beta} \int_{\theta} \sum_z p(\beta, \theta, z, w)}$$

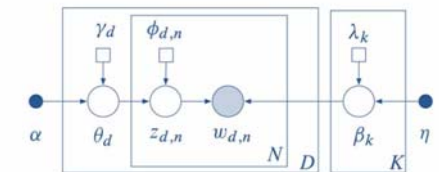
We can't compute the denominator, the marginal $p(w)$, therefore we use **approximate inference**;
However, this do not scale well ...



Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. 2013. Stochastic variational inference. The Journal of Machine Learning Research, 14, (1), 1303-1347.

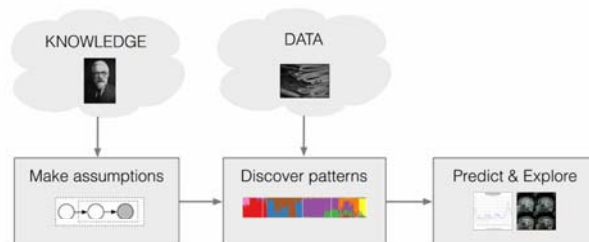


Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. 2013. Stochastic variational inference. The Journal of Machine Learning Research, 14, (1), 1303-1347.

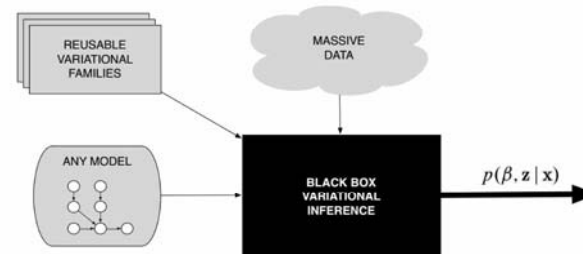


- Sample a document
- Estimate the local variational parameters using the current topics
- Form intermediate topics from those local parameters
- Update topics as a weighted average of intermediate and current topics

Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. 2013. Stochastic variational inference. The Journal of Machine Learning Research, 14, (1), 1303-1347.



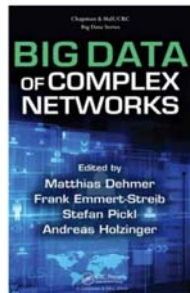
Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. 2013. Stochastic variational inference. The Journal of Machine Learning Research, 14, (1), 1303-1347.



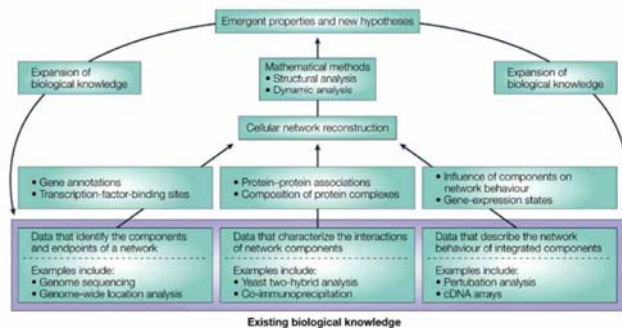
Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. 2013. Stochastic variational inference. The Journal of Machine Learning Research, 14, (1), 1303-1347.

- Flexible and expressive components for building models
- Scalable and generic inference algorithms
- Easy to use software to stretch probabilistic modeling into the health domain
- Topic models are only one approach towards detection of topics in text collections
- More general: Identify re-occurring patterns in data collections generally ...
- Much open work for you in the future ☺

- Particular topic models
 - Stanford topic model toolbox
<http://nlp.stanford.edu/software/tmt>
 - Topic modeling at Princeton
<http://www.cs.princeton.edu/~blei/topicmodeling.html>
 - MALLET (Java) <http://mallet.cs.umass.edu>
 - Network topic models: Bayes-stack
<https://github.com/bgamari/bayes-stack>
 - Gensim (Python) <http://radimrehurek.com/gensim/>
 - R package for Topic models. <http://epub.wu.ac.at/3987/>
- Frameworks for generative models
 - Variational inference: Infer.net
<http://research.microsoft.com/infernet/>
 - Gibbs sampling: OpenBUGS <http://openbugs.net/>



Dehmer, M., Emmert-Streib, F., Pickl, S. & Holzinger, A. (eds.) 2016. Big Data of Complex Networks, Boca Raton, London, New York: CRC Press Taylor & Francis Group.



Nature Reviews | Molecular Cell Biology

Image description find here:
http://www.nature.com/nrm/journal/v6/n2/fig_tab/nrm1570_F1.html

03 Knowledge Representation in Network Medicine

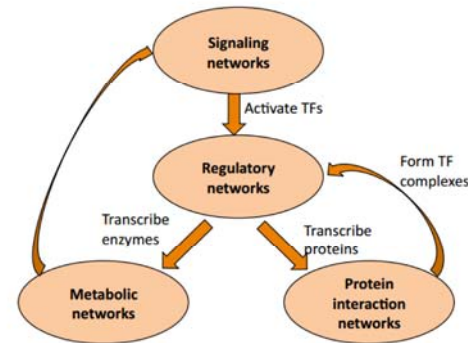


Image credit to Anna Goldenberg, Toronto

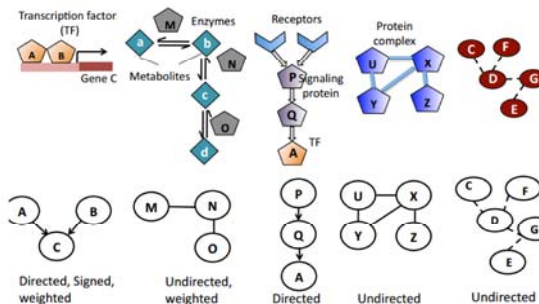


Image credit to Anna Goldenberg, Toronto

Networks = Graphs



<http://www.wired.com/tag/network-science/>

http://www.barabasilab.com/pubs/CCNR-ALB_Publications/200907-24_Science-Decade/200907-24_Science-CoverImage.gif

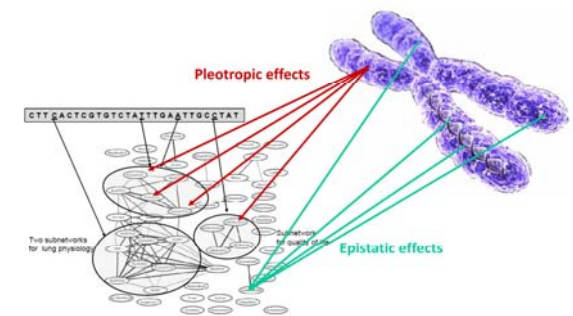
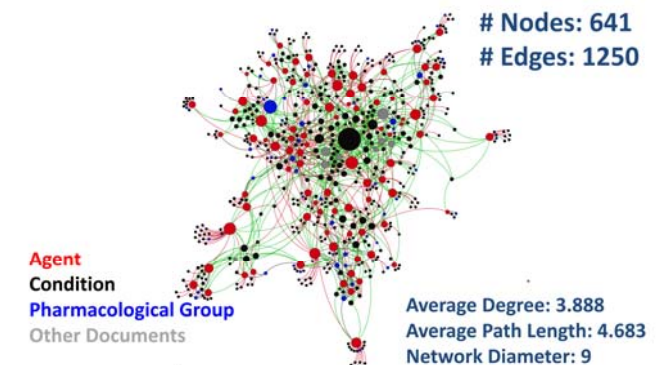


Image credit to Eric Xing, Carnegie Mellon University, Pittsburgh



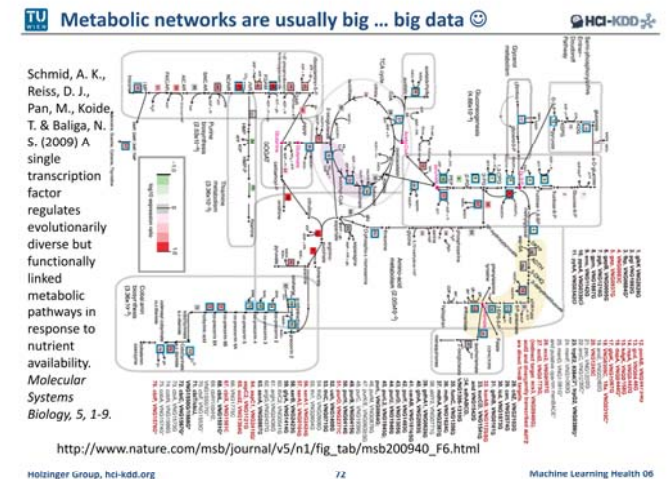
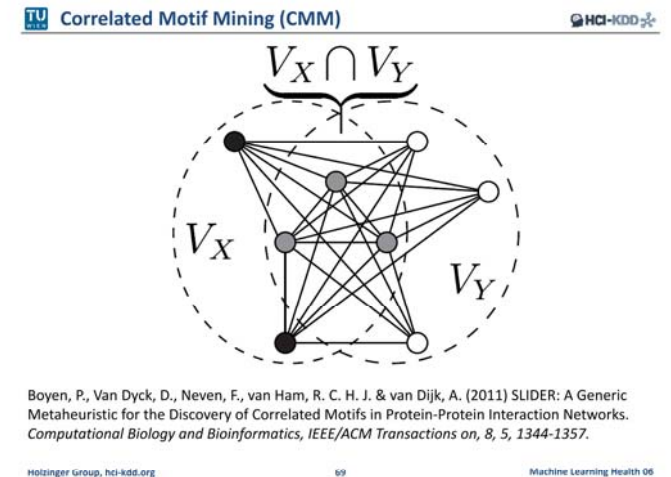
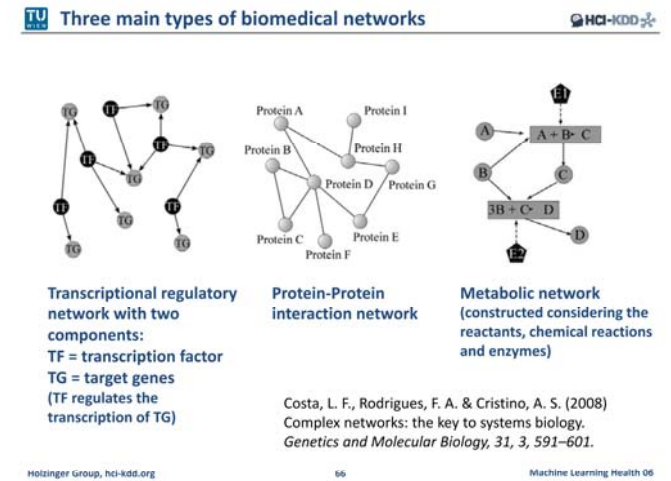
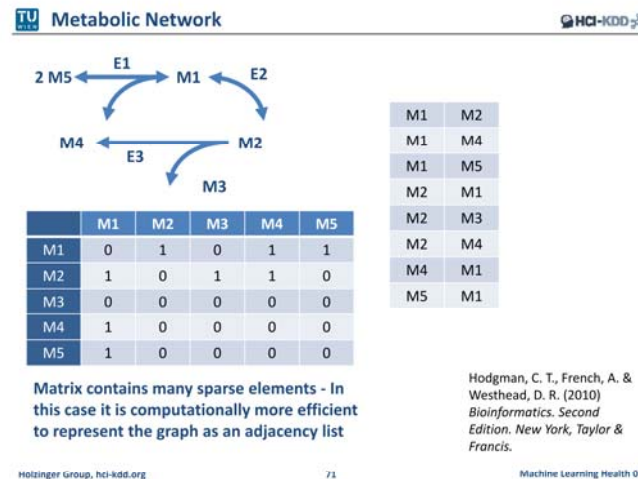
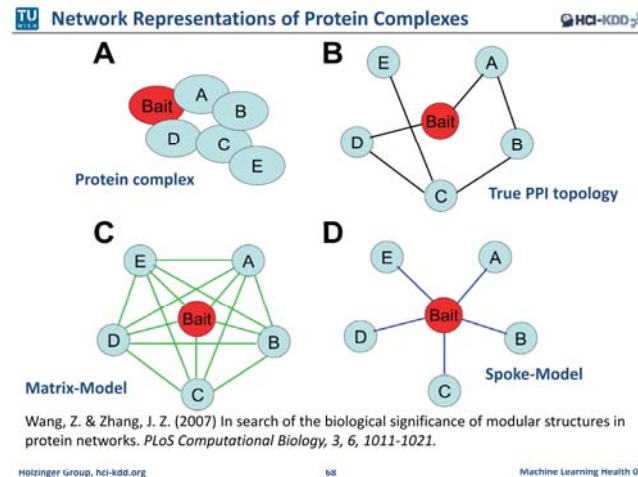
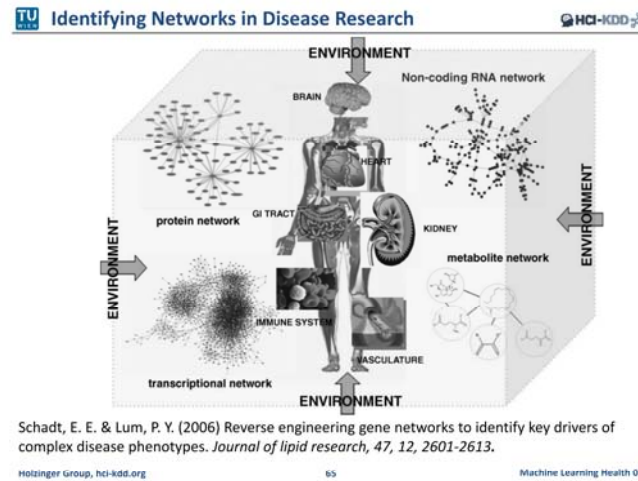
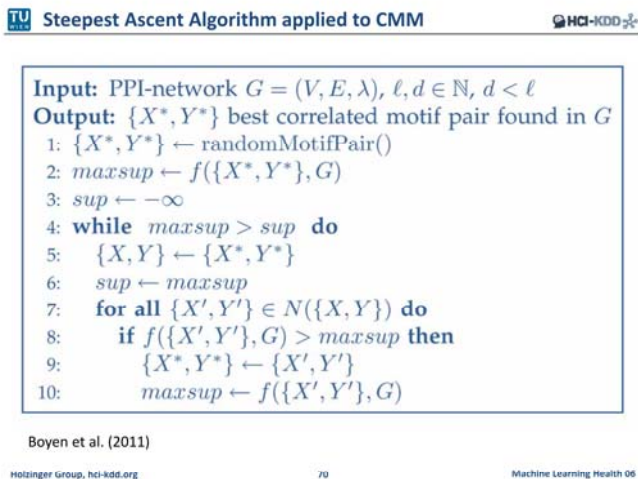
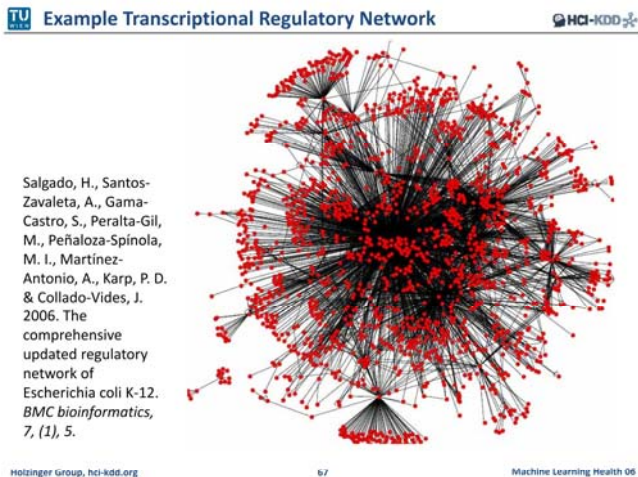
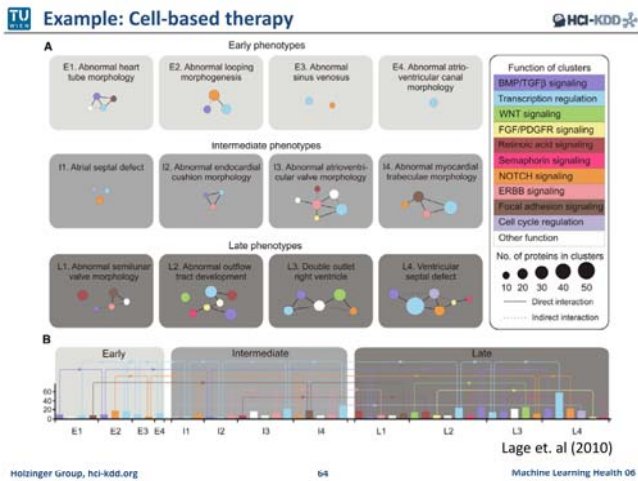
Holzinger, A., Ofner, B., Dehmer, M.: Multi-touch Graph-Based Interaction for Knowledge Discovery on Mobile Devices: State-of-the-Art and Future Challenges. In: LNCS 8401, pp. 241–254, (2014)

55



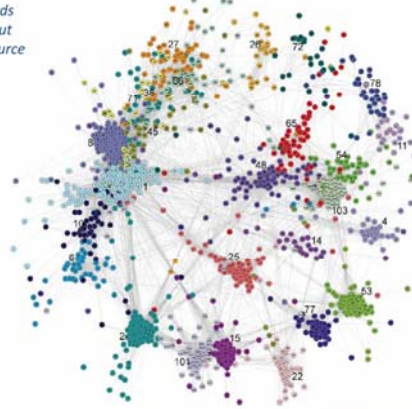
Holzinger Group, hci-kdi





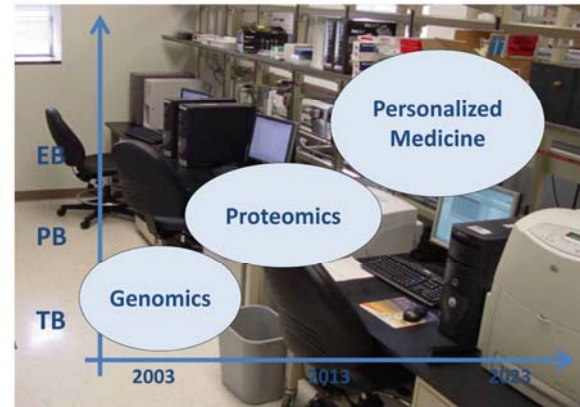
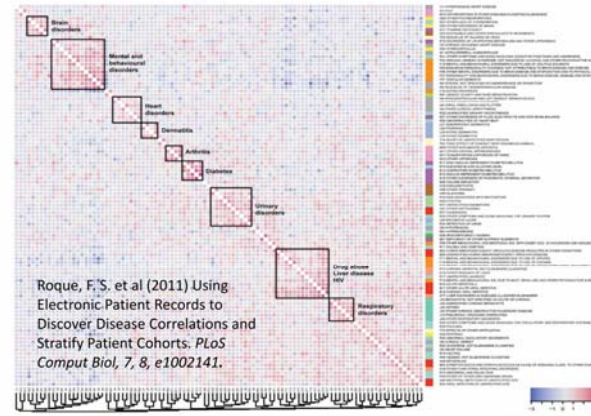
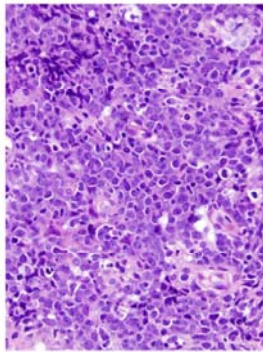
Electronic patient records remain a unexplored, but potentially rich data source for example to discover correlations between diseases.

Roque, F. S., Jensen, P. B., Schmock, H., Dalgaard, M., Andreatta, M., Hansen, T., Sæby, K., Bredkjær, S., Juul, A., Werge, T., Jensen, L. J. & Brunak, S. (2011) Using Electronic Patient Records to Discover Disease Correlations and Stratify Patient Cohorts. *PLoS Computational Biology*, 7, 8, e1002141.

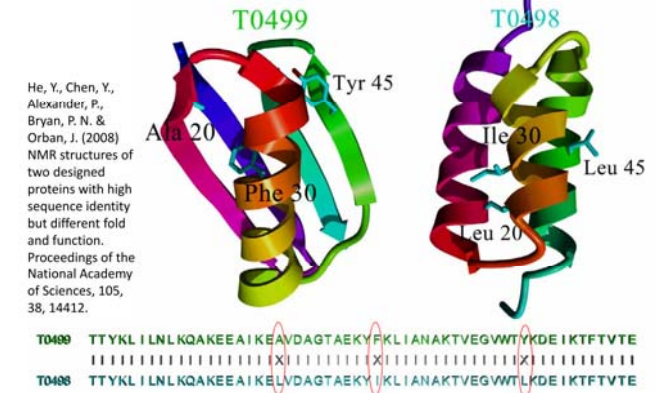
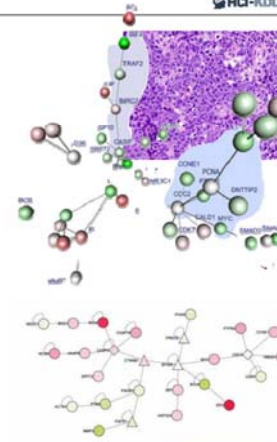
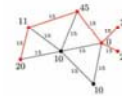


- Homology modeling is a knowledge-based prediction of protein structures.
- In homology modeling a protein sequence with an unknown structure (the target) is aligned with one or more protein sequences with known structures (the templates).
- The method is based on the principle that homologue proteins have similar structures.
- Homology modeling will be extremely important to personalized and molecular medicine in the future.**

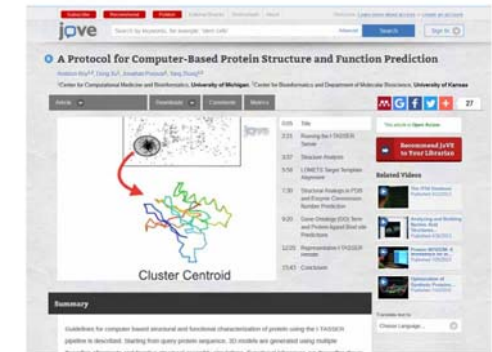
The two main forms of lymphoma are Hodgkin lymphoma and non-Hodgkin lymphoma (NHL). Lymphoma occurs when cells of the immune system called lymphocytes, a type of white blood cell, grow and multiply uncontrollably. Cancerous lymphocytes can travel to many parts of the body, including the lymph nodes, spleen, bone marrow, blood, or other organs, and form a mass called a tumor. The body has two main types of lymphocytes that can develop into lymphomas: B-lymphocytes (B-cells) and T-lymphocytes (T-cells).



- Discover unexplored interactions in PPI-networks and gene regulatory networks
- Learn the structure
- Reconstruct the structure



04 Machine Learning on Graphs Examples



Nodes: proteins
Links: physical interactions (binding)
Puzzling pattern:
Hubs tend to link to small degree nodes.

Why is this puzzling?

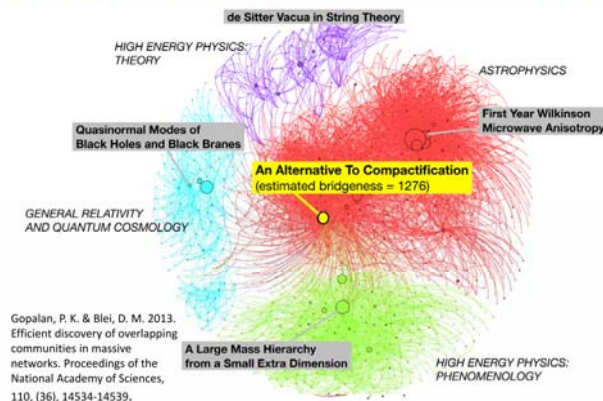
In a random network, the probability that a node with degree k links to a node with degree k' is:

$$p_{kk'} = \frac{kk'}{2L}$$

$$k=50, k'=13, N=1,458, L=1746$$

$$p_{50,13} = 0.15 \quad p_{2,3} = 0.0004$$

Jeong, H., Mason, S. P., Barabasi, A. L. & Oltvai, Z. N. 2001. Lethality and centrality in protein networks. Nature, 411, (6833), 41-42.



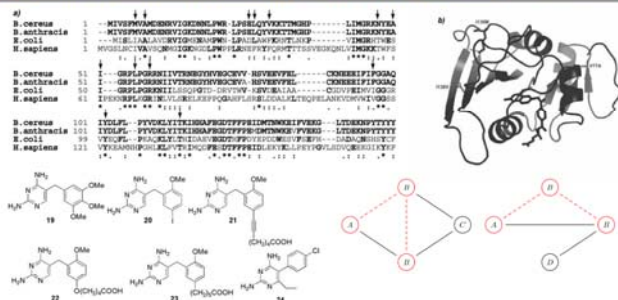
Gopalan, P. K. & Blei, D. M. 2013. Efficient discovery of overlapping communities in massive networks. Proceedings of the National Academy of Sciences, 110, (36), 14534-14539.

- A) Discovery of unexplored interactions
- B) Learning and Predicting the structure
- C) Reconstructing the structure
- Which joint probability distributions does a graphical model represent?
- How can we learn the parameters and structure of a graphical model?



The chemical space

- 10^{60} possible small organic molecules
- 10^{22} stars in the observable universe



How similar are two graphs? How similar is their structure? How similar are their node and edge labels?

Joska, T. M. & Anderson, A. C. 2006. Structure-activity relationships of Bacillus cereus and Bacillus anthracis dihydrofolate reductase: toward the identification of new potent drug leads. Antimicrobial agents and chemotherapy, 50, 3435-3443.

- Similar Property Principle: Molecules having similar structures should have similar activities.
- Structure-based representations: Compare molecules by comparing substructures, e.g.
 - Sets as vectors: Measure similarity by the cosine distance
 - Sets as sets: Measure similarity by the Jaccard distance
 - Sets as points: Measure similarity by Euclidean distance
- Problems: Dimensionality, Non-Euclidean cases

05 Digression: What is similarity?

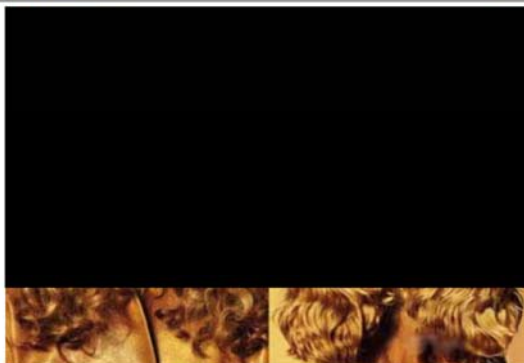
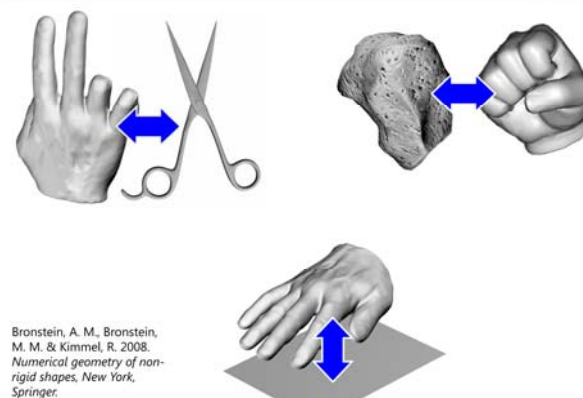
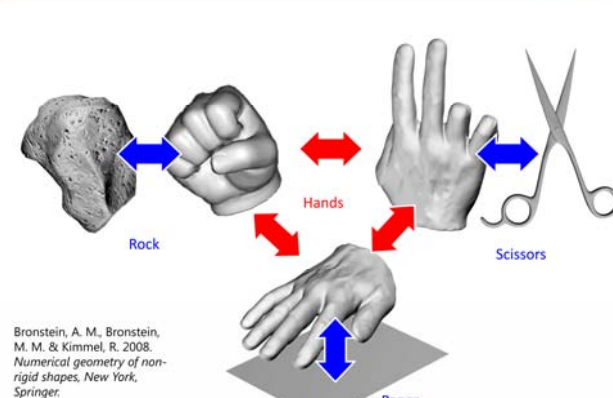


Image credit to Eamonn Keogh (2008)



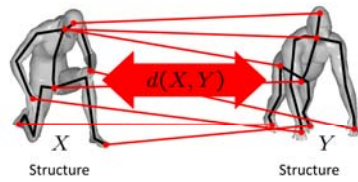
Bronstein, A. M., Bronstein, M. M. & Kimmel, R. 2008. Numerical geometry of non-rigid shapes. New York, Springer.



Bronstein, A. M., Bronstein, M. M. & Kimmel, R. 2008. Numerical geometry of non-rigid shapes. New York, Springer.

Bronstein, A. M., Bronstein, M. M. & Kimmel, R. 2008. *Numerical geometry of non-rigid shapes*, New York, Springer.

<http://www.inf.usi.ch/bronstein/>



Correspondence quality = structure similarity (distortion)

Minimum possible correspondence distortion

Holzinger Group, hci-kdd.org

91

Machine Learning Health 06



Enrico Betti (1823-1892)

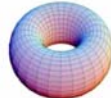
Counts the number of "i-dimensional holes"
bi is the "i-th Betti number"



$b_1=1$
 $b_2=0$



$b_1=0$
 $b_2=1$



$b_1=2$
 $b_2=1$



Emmy Noether (1882-1935)

Betti numbers are computed as dimensions of Boolean vector spaces (E. Noether)
Zomorodian, A. & Carlsson, G. 2005. *Computing Persistent Homology*. *Discrete & Computational Geometry*, 33, (2), 249-274.

Holzinger Group, hci-kdd.org

94

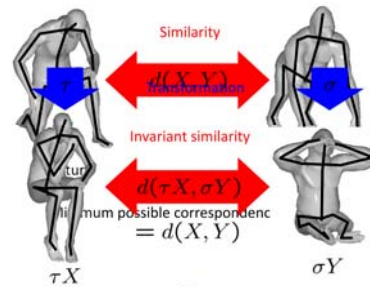
Machine Learning Health 06

06 Review of basic concepts, metrics and measures

Holzinger Group, hci-kdd.org

97

Machine Learning Health 06



Holzinger Group, hci-kdd.org

92

Machine Learning Health 06

- Statement of Vin de Silva (2003), Pomona College:
- Let M be a topological or metric space, known as the *hidden parameter space*;
- let \mathbb{R}^d be a Euclidean space, the *observation space*,
- and let $f: M \rightarrow \mathbb{R}^d$ be a continuous embedding.
- Furthermore, let $X \subset M$ be a finite set of data points, perhaps the realization of a stochastic process, i.e., a family of random variables $\{X_i, i \in I\}$ defined on a probability space (Ω, \mathcal{F}, P) , and denote $Y = f(X) \subset \mathbb{R}^d$ the images of these points under the mapping f .
- We refer to X as *hidden data*, and Y as the *observed data*.
- M, f and X are unknown, but Y is - so can we identify M ?

Holzinger Group, hci-kdd.org

95

Machine Learning Health 06

- In order to understand complex biological systems, the three following key concepts need to be considered:
- (i) **emergence**, the discovery of links between elements of a system because the study of individual elements such as genes, proteins and metabolites is insufficient to explain the behavior of whole systems;
- (ii) **robustness**, biological systems maintain their main functions even under perturbations imposed by the environment; and
- (iii) **modularity**, vertices sharing similar functions are highly connected.
- Network theory can largely be applied for biomedical informatics, because many tools are already available

Holzinger Group, hci-kdd.org

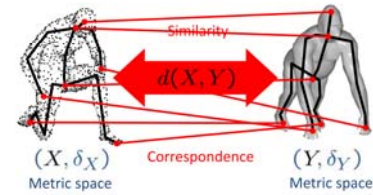
98

Machine Learning Health 06



Mikhail Gromov (1943-)

Gromov, M. (1984) Infinite groups as geometric objects.



Felix Hausdorff (1868-1942)

$$d_{GH}(X, Y) = \frac{1}{2} \min_C \max_{(x_i, y_i) \in C} |\delta_X(x_i, x_j) - \delta_Y(y_i, y_j)|$$

$$\forall x_i \exists y_i \text{ s.t. } (x_i, y_i) \in C \quad \forall y_i \exists x_i \text{ s.t. } (x_i, y_i) \in C$$

Discrete optimization over correspondences is NP hard !

Holzinger Group, hci-kdd.org

93

Machine Learning Health 06



- Mega Problem: To date none of our known methods, algorithms and tools scale to the massive amount and dimensionalities of data we are confronted in practice;
- we need much more research efforts towards making computational topology successful as a general method for data mining and knowledge discovery

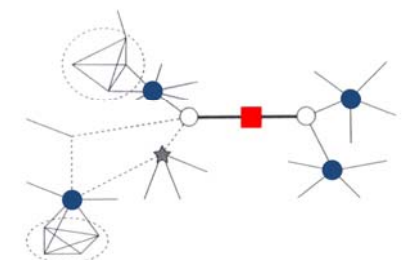
Holzinger, A. 2014. On Topological Data Mining. In: *Lecture Notes in Computer Science*, LNCS 8401. Berlin Heidelberg: Springer, pp. 331-356, doi:10.1007/978-3-662-43968-5_19.

Holzinger Group, hci-kdd.org

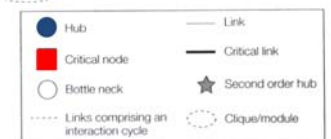
96

Machine Learning Health 06

$G(V, E)$ Graph
 $V \dots$ vertex
 $E \dots$ edge $\{a, b\}$
 $a, b \in V; a \neq b$



Hodgman, C. T., French, A. & Westhead, D. R. (2010) *Bioinformatics*. Second Edition. New York, Taylor & Francis.



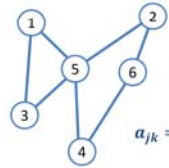
Holzinger Group, hci-kdd.org

99

Machine Learning Health 06

Adjacency (a- 'jā-s'n(t)-sē) Matrix $A = (a_{jk})$

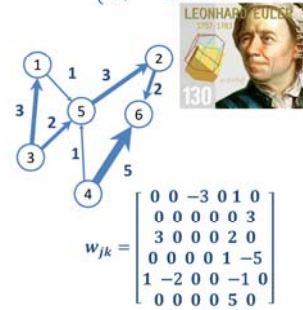
$$a_{jk} = \begin{cases} 1, & \text{if } \{j, k\} \in E \\ 0, & \text{otherwise} \end{cases}$$



$$a_{jk} = \begin{bmatrix} 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$$

Simple graph, symmetric, binary

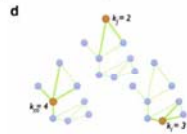
Directed and weighted



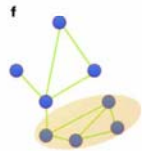
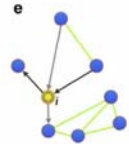
$$w_{jk} = \begin{bmatrix} 0 & 0 & -3 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \\ 3 & 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 1 & -5 \\ 1 & -2 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 5 & 0 \end{bmatrix}$$

For more information: Diestel, R. (2010) *Graph Theory, 4th Edition. Berlin, Heidelberg, Springer.*

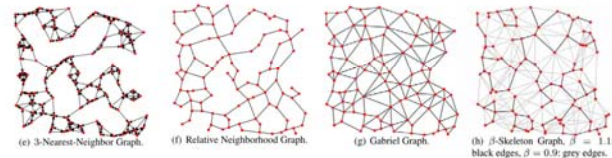
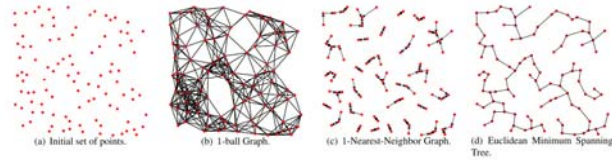
- Centrality (d) = the level of "betweenness- centrality" of a node i ("hub-node in Slide 28);



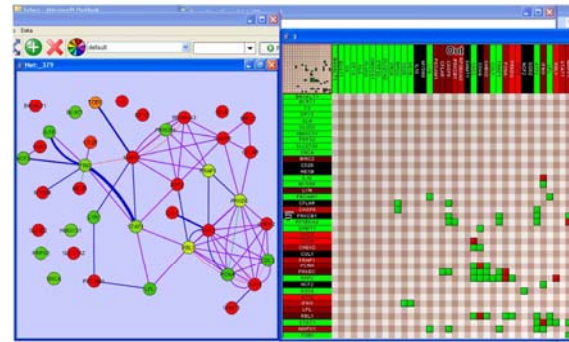
- Nodal degree (e) = number of links connecting i to its neighbors: $k_i = \sum_j a_{ij}$



Modularity (f) = describes the possible formation of communities in the network, indicating how strong groups of nodes form relative isolated sub-networks within the full network (refer also to Slide 5-8).

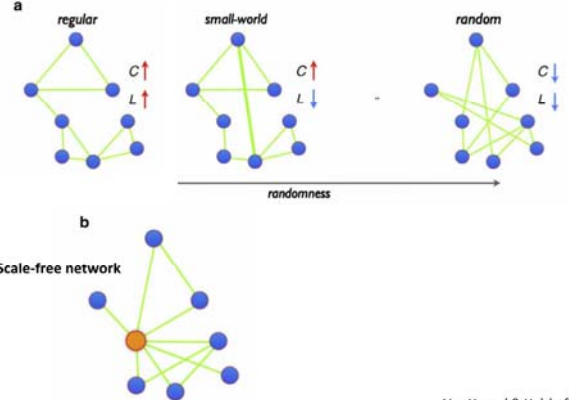


Lézoray, O. & Grady, L. 2012. Graph theory concepts and definitions used in image processing and analysis. In: Lézoray, O. & Grady, L. (eds.) *Image Processing and Analysing With Graphs: Theory and Practice. Boca Raton (FL): CRC Press*, pp. 1-24.



Jean-Daniel Fekete http://wiki.cytoscape.org/InfoVis_Toolkit

Fekete, J.-D. The Infovis toolkit. Information Visualization, INFOVIS 2004, 2004. IEEE, 167-174.

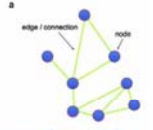


Van Heuvel & Hulshoff (2010)

07 How do you get point cloud data from natural images?

Order = total number of nodes n ; Size = total number of links (a):

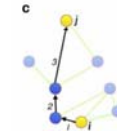
$$\sum_i \sum_j a_{ij}$$



Clustering Coefficient (b) = the degree of concentration of the connections of the node's neighbors in a graph and gives a measure of local inhomogeneity of the link density:

$$C_i = \frac{2t_i}{k_i(k_i - 1)}$$

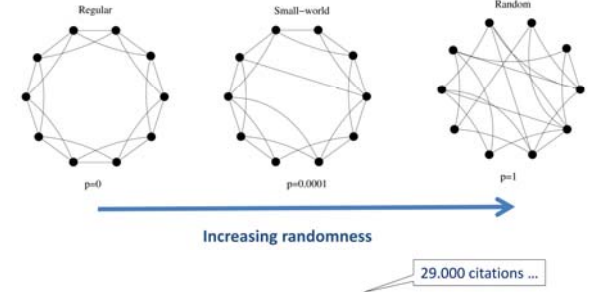
$$C = \frac{1}{n} \sum_i C_i$$



Path length (c) = is the arithmetical mean of all the distances:

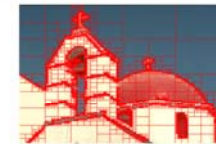
$$l = \frac{1}{n(n-1)} \sum_{i,j} d_{ij}$$

Costa, L. F., Rodrigues, F. A., Traviesso, G. & Boas, P. R. V. (2007) Characterization of complex networks: A survey of measurements. *Advances in Physics*, 56, 1, 167-242.



Watts, D. J. & Strogatz, S. (1998) Collective dynamics of small-world networks. *Nature*, 393, 6684, 440-442.

Milgram, S. 1967. The small world problem. *Psychology today*, 2, (1), 60-67.



a) quadtree tessellation



b) RAG assoc. to the quadtree



c) Watershed Algorithm



d) SLIC superpixels

Lézoray, O. & Grady, L. 2012. Graph theory concepts and definitions used in image processing and analysis. In: Lézoray, O. & Grady, L. (eds.) *Image Processing and Analysing With Graphs: Theory and Practice. Boca Raton (FL): CRC Press*, pp. 1-24.

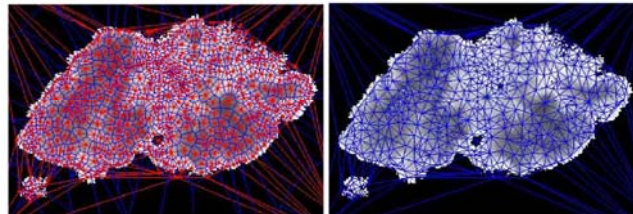
Algorithm 4.2 Watershed transform w.r.t. topographical distance based on image integration via the Dijkstra-Moore shortest paths algorithm.

```

1: procedure ShortestPathWatershed;
2: INPUT: lower complete digital grey scale image  $G = (V, E, im)$  with cost function cost.
3: OUTPUT: labelled image  $lab$  on  $V$ .
4: #define WSHED 0
5: (* Uses distance image  $dist$ . On output,  $dist[v] = im[v]$ , for all  $v \in V$  *)
6:
7: for all  $v \in V$  do (* Initialize *)
8:    $lab[v] \leftarrow 0$ ;  $dist[v] \leftarrow \infty$ 
9: end for
10: for all local minima  $m_i$  do
11:   for all  $v \in m_i$  do
12:      $lab[v] \leftarrow i$ ;  $dist[v] \leftarrow im[v]$  (* initialize distance with values of minima *)
13:   end for
14: end for
15: while  $V \neq \emptyset$  do
16:    $u \leftarrow \text{GetMinDist}(V)$  (* find  $u \in V$  with smallest distance value  $dist[u]$  *)
17:    $V \leftarrow V \setminus \{u\}$ 
18:   for all  $v \in V$  with  $(u, v) \in E$  do
19:     if  $dist[u] + cost[u, v] < dist[v]$  then
20:        $dist[v] \leftarrow dist[u] + cost[u, v]$ 
21:        $lab[v] \leftarrow lab[u]$ 
22:     else if  $lab[v] \neq \text{WSHED}$  and  $dist[u] + cost[u, v] = dist[v]$  and  $lab[v] \neq lab[u]$  then
23:        $lab[v] = \text{WSHED}$ 
24:     end if
25:   end for
26: end while

```

Meijster, A. & Roerdink, J. B. A proposal for the implementation of a parallel watershed algorithm. Computer Analysis of Images and Patterns, 1995. Springer, 790-795.



Holzinger, A., Malle, B. & Giuliani, N. 2014. On Graph Extraction from Image Data. In: Slezak, D., Peters, J. F., Tan, A.-H. & Schwabe, L. (eds.) Brain Informatics and Health, BIH 2014. Lecture Notes in Artificial Intelligence, LNAI 8609. Heidelberg, Berlin: Springer, pp. 552-563.

For Voronoi please refer to: Aurenhammer, F. 1991. Voronoi Diagrams - A Survey of a fundamental geometric data structure. *Computing Surveys*, 23, (3), 345-405.

For Delaunay please refer to: Lee, D.-T. & Schachter, B. J. 1980. Two algorithms for constructing a Delaunay triangulation. *Intl. Journal of Computer & Information Sciences*, 9, (3), 219-242.

- 1) Transformation into a topographic map
 - Convert gray values into height information
- 2) Finding local minima
 - Inspecting small regions in sequence
- 3) Finding catchment basins
 - Algorithm simulating flooding
 - Graph algorithms such as Minimum Spanning Trees
- 4) Erecting watersheds
 - Artificial divide between catchment basins
 - Final segmentation lines

- More expressive data structures
- Find novel connections between data objects
- Fit for applying graph based machine learning techniques
- New approaches (Belief Propagation, global understanding from local properties)

Bunke, H.: Graph-based tools for data mining and machine learning. In Perner, P., Rosenfeld, A., eds.: Machine Learning and Data Mining in Pattern Recognition, Proceedings. Volume 2734 of Lecture Notes in Artificial Intelligence. Springer-Verlag Berlin, (Berlin) 7-19

Holzinger, A., Blanchard, D., Bloice, M., Holzinger, K., Palade, V., Rabadan, R.: Darwin, lamarck, or baldwin: Applying evolutionary algorithms to machine learning techniques. In: The 2014 IEEE/WIC/ACM International Conference on Web Intelligence (WI 2014), IEEE (2014) in print

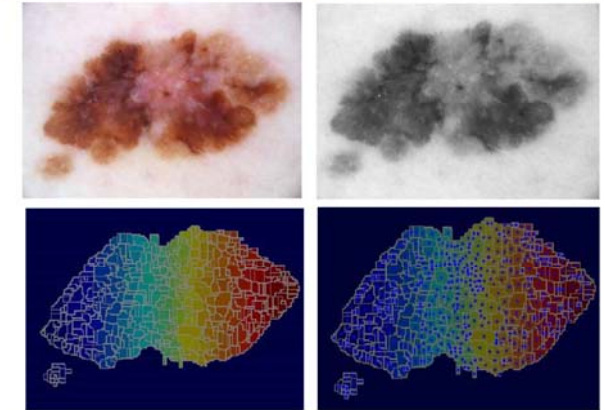
7	4	8	12	11	3	→	m	←	←	→	m	0	0	0	0	1	1
7	7	8	12	11	7	↗	↑	↖	←	↗	↑	0	0	0	0	1	1
13	13	15	16	16	13	↑	↑	↖	↖	↗	↑	0	0	0	0	1	1
19	19	18	17	15	7	↑	↑	↑	→	↘	↓	0	0	0	2	2	2
20	18	17	16	15	5	→	→	→	→	→	m	2	2	2	2	2	2

(a) The original image

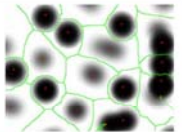
(b) Each pixel connect to lowest minimum

(c) The Image with labels

Connects each pixel to the lowest neighbor pixel, all pixel connected to same lowest neighbor pixel form a segment



- Topographic maps => landscapes with height structures
- Segmentation into regions of pixels
- Assuming drops of water raining on the map
- Following paths of descent
- Lakes called catchment basins
- Also possible: Flooding based
- Needs Topographical distance measures (MST)



Vincent, L. & Soille, P. 1991. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. *IEEE transactions on pattern analysis and machine intelligence*, 13, (6), 583-598.

- Region Merging
 - Based on Kruskals MST algorithm
 - Takes input image as natural graph with vertices := pixels and edges := pixel neighborhoods
 - Visits edges in ascending order of weight and merges regions if they satisfy a certain criterion
 - Flexible as merging criterion can be adapted as desired (for amount, size, or shape of resulting regions)

Felzenszwalb, P.F., Huttenlocher, D.P.: Efficient graph-based image segmentation. *International Journal of Computer Vision* 59 (2004) 167-181

- We want to find “interesting” novel patterns (rules, anomalies, outliers, similarities, ...)
- Problem #1: How to get a graph?
- Problem #2: How do graphs evolve?
- Problem #3: What tools to apply?
- Problem #4: Scalability to TB, PB, EB ...
- **Success is in repeatability and scalability**

Questions

Sample Questions (2/3)

- Briefly describe the stochastic variational inference algorithms!
- What is the principle of a bandit?
- How does a multi-armed bandit (MAB) work?
- In which ways can a MAB represent knowledge?
- What is the main problem of a clinical trial – and maybe the main problem in clinical medicine?
- Why are rare diseases both important and relevant? Describe an example disease!
- What is the big problem in clinical trials for rare diseases?
- What did Richard Bellman (1956) describe with dynamic programming?
- Why are graph bandits a hot topic for ML research?

- Study of complex networks started in the 1990s with the insight that real networks contain properties not present in random (Erdős-Renyi) networks.
- Meanwhile networks and network-based approaches form an integral part of many studies throughout the sciences.
- Graph-Theory provides powerful tools to organize data structurally and in combination with statistical and machine learning methods allows a meaningful analysis of underlying processes.
- For instance, a mapping of causal disease genes and disorders as made available by the OMIM database provided novel insights into disease patterns, as recently demonstrated by investigating the diseasome (<http://diseasome.eu>).

Sample Questions (1/3)

- Describe the clinical decision making process!
- Which type of graph is particularly useful for inference and learning?
- What is the key challenge in the application of graphical models for health informatics?
- What was Judea Pearl (1988) discussing in his paper, for which he received the Turing award?
- What main difficulties arise during breast cancer prognosis?
- What can be done to increase the robustness of prognostic cancer tests?
- Inference in Bayes Nets is NP-complete, but there are certain cases where it is tractable, which ones?

Solutions of the Quiz

- 1= this is a factor graph of an undirected graph – we have seen this in protein networks (refer to slide Nr. 70 in lecture 5). Factor graph is bipartite and has two types of nodes: Variables, which can be either evidence variables (when we know its value) or query variables (when the value is unknown and we want to predict the value); and factors, which define the relationship between variables in the graph. Each factor can be connected to many variables and comes with a factor function to define the relationship between these variables. For example, if a factor node is connected to two variables nodes A and B, a possible factor function could be $\text{imply}(A,B)$, meaning that if the random variable A takes value 1, then so must the random variable B. Each factor function has a weight associated with it, which describes how much influence the factor has on its variables in relative terms. For more information please consult: <http://deepdive.stanford.edu/inference>
- 2= this is the decomposition of a tree, rooted at nodes into subtrees
- 3= an example for machine translation, Image credit to Kevin Gimpel, Carnegie Mellon University
- 4= the famous expectation-utility theory according to von Neumann and Morgenstern (1954): a decision-maker faced with risky (probabilistic) outcomes of different choices will behave as if he is maximizing the expected value of some function defined over the potential outcomes at some specified point in the future.
- 5= MYCIN – expert system that used early AI (rule-based) to identify bacteria causing severe infections, such as bacteremia and meningitis, and to recommend antibiotics, with the dosage adjusted for patient's body weight – the name derived from the antibiotics themselves, as many antibiotics have the suffix “-mycin”.
- 6= metabolic and physical processes that determine the physiological and biochemical properties of a cell. These networks comprise the chemical reactions of metabolism, the metabolic pathways, as well as the regulatory interactions that guide these reactions.
- 7= With the sequencing of complete genomes, it is now possible to reconstruct the network of biochemical reactions in many organisms, from bacteria to human. Several of these networks are available online, e.g. Kyoto Encyclopedia of Genes and Genomes (KEGG), EcoCyc, BioCyc etc. Metabolic networks are powerful tools for studying and modelling metabolism.



Thank you!

Sample Questions (2/3)

- Why do we want to apply ML to graphs?
- Describe typical ML tasks on the example of blood cancer cells!
- If you have a set of points – which similarity measures are useful?
- Why is graph comparison in the medical domain useful?
- Why is the Gromov-Hausdorff distance useful?
- What is the central goal of a generative probabilistic model?
- Describe the LDA-model and its application for topic modelling!

Appendix

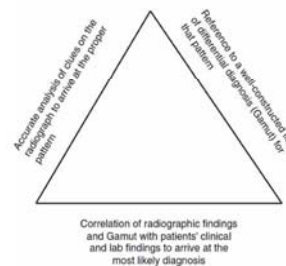
1) Reasoning under Uncertainty

MYCIN – rule based system - certainty factors

- MYCIN is a rule-based Expert System, which is used for therapy planning for patients with bacterial infections
- Goal oriented strategy (“Rückwärtsverkettung”)
- To every rule and every entry a certainty factor (CF) is assigned, which is between 0 und 1
- Two measures are derived:
 - MB: measure of belief
 - MD: measure of disbelief
- Certainty factor – CF of an element is calculated by:

$$CF[h] = MB[h] - MD[h]$$
- CF is positive, if more evidence is given for a hypothesis, otherwise CF is negative
- CF[h] = +1 -> h is 100 % true
- CF[h] = -1 -> h is 100% false

Gamuts: Triangulation to find diagnoses



Reeder, M. M. & Felson, B. 2003. *Reeder and Felson's gamuts in radiology: comprehensive lists of roentgen differential diagnosis*, New York, Springer Verlag.

Gamut F-137 PHRENIC NERVE PARALYSIS OR DYSFUNCTION

COMMON

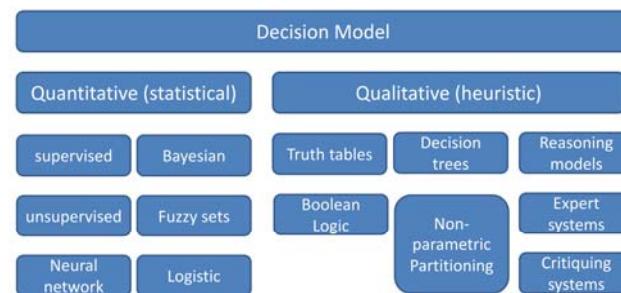
1. Iatrogenic (eg, surgical injury; chest tube; therapeutic avulsion or injection; subclavian vein puncture)
2. Infection (eg, tuberculosis; fungus disease; abscess)
3. Neoplastic invasion or compression (esp. carcinoma of lung)

UNCOMMON

1. Anarhythm, aortic or other
2. Birth trauma (Erb's palsy)
3. Herpes zoster
4. Neuritis, peripheral (eg, diabetic neuropathy)
5. Neurologic disease, (eg, hemiplegia; encephalitis; polio; Guillain-Barré S.)
6. Pneumonia
7. Trauma

Reference

1. Prasad S, Adhaya BH. Transient paralysis of the phrenic nerve associated with head injury. JAMA 1976;236:2532-2533



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

Original Example from MYCIN

h_1 = The identity of ORGANISM-1 is streptococcus
 h_2 = PATIENT-1 is febrile
 h_3 = The name of PATIENT-1 is John Jones

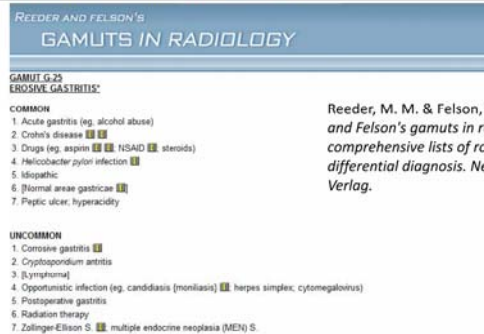
$CF[h_1, E] = .8$: There is strongly suggestive evidence (.8) that the identity of ORGANISM-1 is streptococcus

$CF[h_2, E] = -.3$: There is weakly suggestive evidence (.3) that PATIENT-1 is not febrile

$CF[h_3, E] = +1$: It is definite (1) that the name of PATIENT-1 is John Jones

Shortliffe, E. H. & Buchanan, B. G. (1984) *Rule-based expert systems: the MYCIN experiments of the Stanford Heuristic Programming Project*. Addison-Wesley.

Example - Gamuts in Radiology



Reeder, M. M. & Felson, B. (2003) *Reeder and Felson's gamuts in radiology: comprehensive lists of roentgen differential diagnosis*. New York, Springer Verlag.

* Superficial erosions or aphthoid ulcerations seen especially with double contrast technique.

[] This condition does not actually cause the gamuted imaging finding, but can produce imaging changes that simulate it.

<http://rfs.acr.org/gamuts/data/G-25.htm>

- The information available to humans is often imperfect – imprecise - uncertain.
- This is especially in the medical domain the case.
- An **human agent** can cope with deficiencies.
- Classical logic permits only **exact reasoning**:
- IF A is true THEN A is non-false and IF B is false THEN B is non-true
- Most real-world problems do not provide this exact information, mostly it is inexact, incomplete, uncertain and/or **un-measurable!**

MYCIN was *no* success in the clinical practice

<https://www.youtube.com/watch?v=IVGWM0CKNWA> (“real nurse triage”)



Reasoning under uncertainty

- Take patient information, e.g., observations, symptoms, test results, -omics data, etc. etc.
- Reach conclusions, and predict into the future, e.g. how likely will the patient be re-admissioned
- Prior = belief before making a particular observation
- Posterior = belief after making the observation and is the prior for the next observation – intrinsically incremental

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

- Holzinger Group,
- hci-kdd.org

Holzinger Group, hci-kdd.org

1

