

Andreas Holzinger
185.A83 Machine Learning for Health Informatics
2019S, VU, 2.0 h, 3.0 ECTS
Lecture 04 – Dienstag, 02.04.2019



From Decision Making under uncertainty to graphical models and MCMC

andreas.holzinger AT tuwien.ac.at

<https://hci-kdd.org/machine-learning-for-health-informatics-class-2019>



Why Explainability? Why Causability?

- 00 Reflection from last lecture
- 01 Decision Making under uncertainty
- 02 Some Basics of Markov Processes
- 03 Some Basics of Concept Learning
- 04 Some Basics of Graphs/Networks and Challenges
- 05 Bayes Nets
- 06 Probabilistic Programming
- 07 Markov Chain Monte Carlo (MCMC)
- 08 Metropolis Hastings Algorithm

01 Reflection

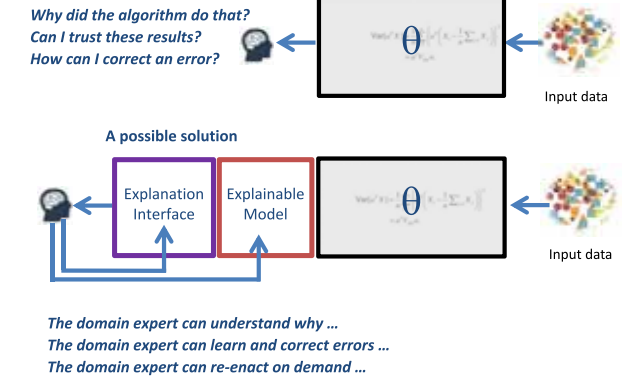


**Causability :=
a property of a person
(Human Intelligence)**

**Explainability :=
a property of a system
(Artificial Intelligence)**

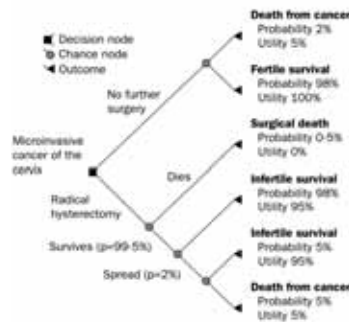
<https://onlinelibrary.wiley.com/doi/full/10.1002/widm.1312>

Andreas Holzinger et al. 2019. Causability and Explainability of AI in Medicine. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 10(1):1002, widm.1312.



- 1) Gradients
- 2) Sensitivity Analysis
- 3) Decomposition Relevance Propagation
 - Pixel-RP, Layer-RP, Deep Taylor Decomposition, ...
- 4) Optimization
 - Local Interpretable Model-Agnostic Explanations (LIME)
 - Black Box Explanations tr. Transparent Approximations (BETA)
- 5) Deconvolution and Guided Backpropagation
- 6) Concept Activation Vectors CAV

Andreas Holzinger LV 706.315 From explainable AI to Causability, 3 ECTS course at Graz University of Technology
<https://hci-kdd.org/explainable-ai-causability-2019> (course given since 2016)



Physician treating a patient
approx. 480 B.C.
Beazley (1963), Attic Red-figured
Vase-Painters, 813, 96.
Department of Greek, Etruscan
and Roman Antiquities, Sully, 1st
floor, Campana Gallery, room 43
Louvre, Paris

Elwyn, G., Edwards, A., Eccles, M. & Rovner, D. 2001. Decision analysis in patient care.
The Lancet, 358, (9281), 571-574.



01 Decision Making under uncertainty

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



... permanent decision making under uncertainty!

$d \dots$ data

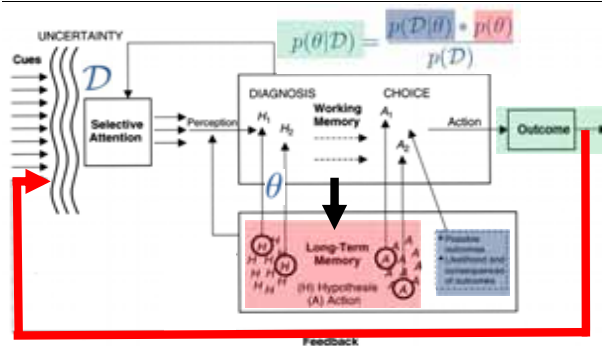
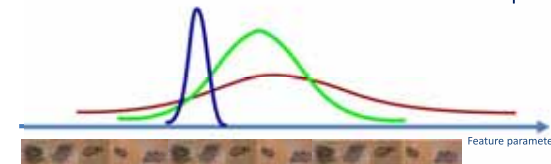
$h \dots$ hypotheses

$\mathcal{H} \dots \{H_1, H_2, \dots, H_n\} \quad \forall h, d \dots$

$$p(h|d) = \frac{\text{Likelihood } p(d|h) * \text{Prior Probability } p(h)}{\sum_{h \in \mathcal{H}} p(d|h) p(h)}$$

Posterior Probability

Problem in $\mathbb{R}^n \rightarrow$ complex



Wickens, C. D. (1984) *Engineering psychology and human performance*. Columbus (OH), Charles Merrill, modified by Holzinger, A.

You are talking to you colleague and want to refer to the middle object – which wording would you prefer: circle or blue?



Frank, M. C. & Goodman, N. D. 2012. Predicting pragmatic reasoning in language games. *Science*, 336, (6084), 998-998, doi:10.1126/science.1218633.

```
var literalListener = function(property){
  Infer(function(){
    var object = refPrior(context)
    condition(object[property])
    return object
  })
}

var speaker = function(object) {
  Infer(function(){
    var property = propPrior()
    condition(
      utterance ==
      object[property]
    )
  })
}

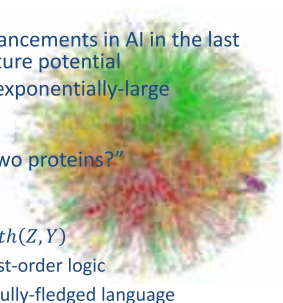
var listener = function(property) {
  Infer(function(){
    var object = refPrior(context)
    condition(utterance ==
      sample(speaker(object)))
    return object
  })
}
```

Goodman, N. D. & Frank, M. C. 2016. Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20, (11), 818-829.

- PGM can be seen as a combination between
- Graph Theory + Probability Theory + Machine Learning**

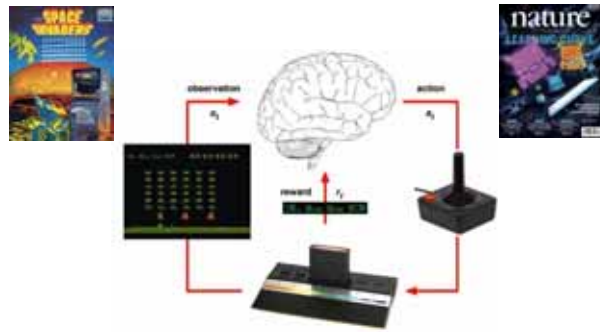
- One of the most exciting advancements in AI in the last decades – with enormous future potential
- Compact representation for exponentially-large probability distributions
- Example Question: "Is there a path connecting two proteins?"

- $Path(X, Y) := edge(X, Y)$
- $Path(X, Y) := edge(X, Y), path(Z, Y)$
- This can NOT be expressed in first-order logic
- Would need a Turing-complete fully-fledged language



2) Some basics of Markov Processes in Machine Learning

- Markov processes are ...
- random processes in which the future, given the present, is independent of the past!
- one of the most important classes of random processes!

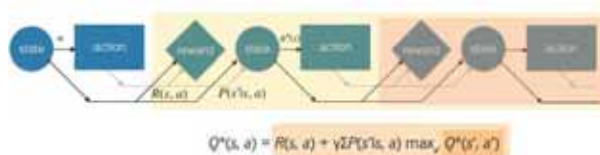


Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S. & Hassabis, D. 2015. Human-level control through deep reinforcement learning. Nature, 518, (7540), 529-533. doi:10.1038/nature14236

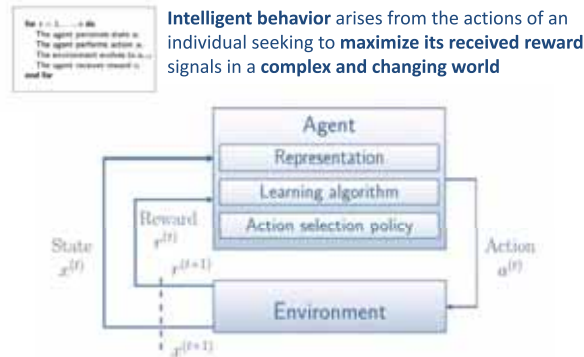


Kaelbling, L. P., Littman, M. L. & Moore, A. W. 1996. Reinforcement learning: A survey. Journal of Artificial Intelligence Research, 4, 237-285.

- Markov decision processes specify setting and tasks
- Planning methods use knowledge of P and R to compute a good policy π
- Markov decision process model captures both sequential feedback and the more specific one-shot feedback (when $P(s'|s, a)$ is independent of both s and a)



Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. Nature, 521, (7553), 445-451.



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

- 1) Oversees
- 2) Executes
- 3) Receives Reward
- Executes action A_t :
- $O_t = s a_t = s_t$
- Agent state = environment state = information state
- Markov decision process (MDP)

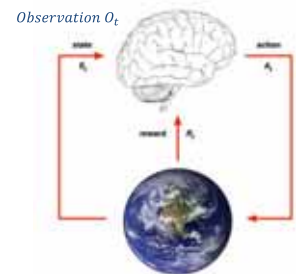


Image credit to David Silver, UCL

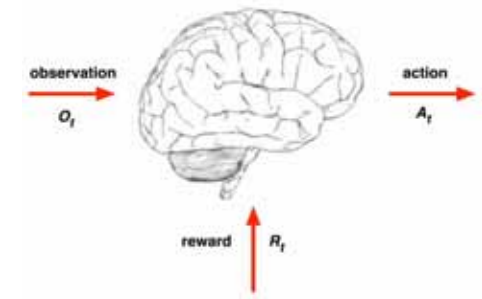
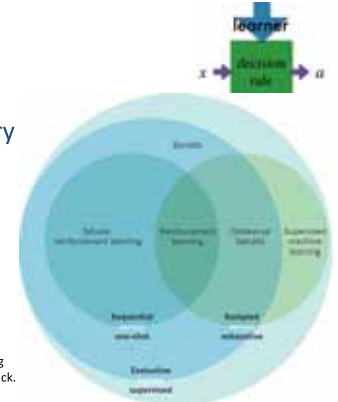


Image credit to David Silver, UCL

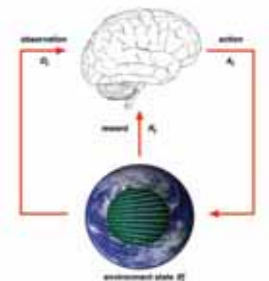
- Supervised: Learner told best a
- Exhaustive: Learner shown every possible x
- One-shot: Current x independent of past a



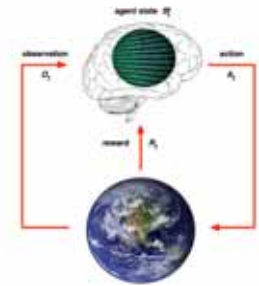
Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. Nature, 521, (7553), 445-451.

- i.e. whatever data the environment uses to pick the next observation/reward
- The environment state is not usually visible to the agent
- Even if S is visible, it may contain irrelevant information
- A State S_t is Markov iff:

$$\mathbb{P}[S_{t+1}|S_t] = \mathbb{P}[S_{t+1}|S_1, \dots, S_t]$$



- i.e. whatever information the agent uses to pick the next action
- it is the information used by reinforcement learning algorithms
- It can be any function of history:
- $S = f(H)$



$$H_t = O_1, R_1, A_1, \dots, A_{t-1}, O_t, R_t$$

- RL agent components:
 - Policy: agent's behaviour function
 - Value function: how good is each state and/or action
 - Model: agent's representation of the environment
- Policy as the agent's behaviour
 - is a map from state to action, e.g.
 - Deterministic policy: $a = (s)$
 - Stochastic policy: $(a|s) = P[A_t = a|S_t = s]$
- Value function is prediction of future reward:

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

- Partial observability: when agent only indirectly observes environment (robot which is not aware of its current location; good example: Poker play: only public cards are observable for the agent):
- Formally this is a Partially Observable Markov Decision Process (POMDP):
 - Agent must construct its own state representation S_t , for example:

- Complete history: $S_t^a = H_t$
- Beliefs of environment state: $S_t^b = (P[S_t^e = s^1], \dots, P[S_t^e = s^n])$
- Recurrent neural network: $S_t^c = \sigma(S_{t-1}^c W_s + O_t W_o)$

3) Some basics of Concept Learning

- **Deductive Reasoning** = Hypothesis > Observations > Logical Conclusions (general \rightarrow specific – proven correctness)
 - DANGER: Hypothesis must be correct! DR defines whether the truth of a conclusion can be determined for that rule, based on the truth of premises: $A=B, B=C$, therefore $A=C$
- **Inductive reasoning** = makes broad generalizations from specific observations (specific \rightarrow general – not proven correctness)
 - DANGER: allows a conclusion to be false if the premises are true
 - generate hypotheses and use DR for answering specific questions
- **Abductive reasoning** = inference = to get the best explanation from an incomplete set of preconditions.
 - Given a true conclusion and a rule, it attempts to select some possible premises that, if true also, may support the conclusion, though not uniquely.
 - Example: "When it rains, the grass gets wet. The grass is wet. Therefore, it might have rained." This kind of reasoning can be used to develop a hypothesis, which in turn can be tested by additional reasoning or data.

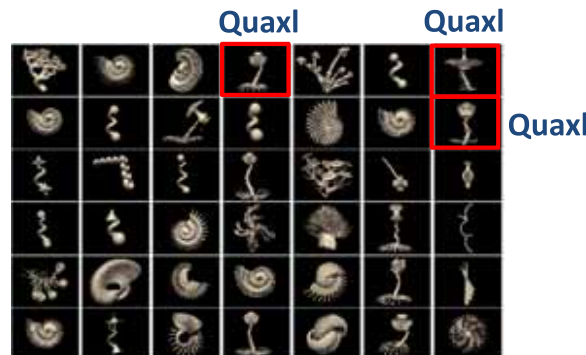
- Bruner, Goodnow, and Austin (1956) published "A Study of Thinking", which became a landmark in cognitive science and has much influence on machine learning.
 - Rule-Based Categories
 - A concept specifies conditions for membership



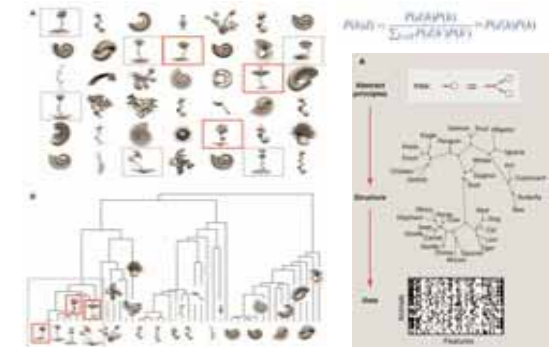
Jerome S. Bruner, Jacqueline J. Goodnow & George A. Austin 1986. A Study of Thinking, Transaction Books.



Salakhutdinov, R., Tenenbaum, J. & Torralba, A. 2012. One-shot learning with a hierarchical nonparametric Bayesian model. Journal of Machine Learning Research, 27, 195-207.



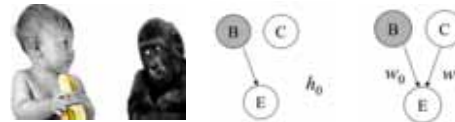
Salakhutdinov, R., Tenenbaum, J. & Torralba, A. 2012. One-shot learning with a hierarchical nonparametric Bayesian model. Journal of Machine Learning Research, 27, 195-207.



Tenenbaum, J. B., Kemp, C., Griffiths, T. L. & Goodman, N. D. 2011. How to grow a mind: Statistics, structure, and abstraction. Science, 331, (6022), 1279-1285.

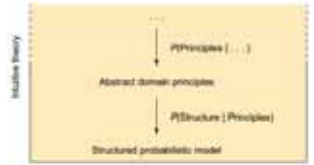
- which is highly relevant for ML research, concerns the factors that determine the subjective difficulty of concepts:
- Why are some concepts psychologically extremely simple and easy to learn,
- while others seem to be extremely difficult, complex, or even incoherent?
- These questions have been studied since the 1960s but are still unanswered ...

Feldman, J. 2000. Minimization of Boolean complexity in human concept learning. Nature, 407, (6804), 630-633, doi:10.1038/35036586.

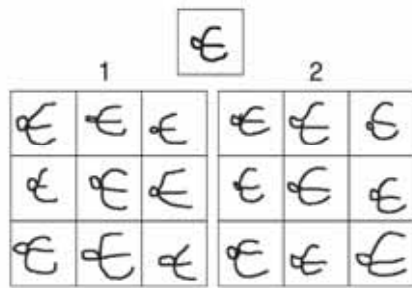


- Cognition as probabilistic inference
 - Visual perception, language acquisition, motor learning, associative learning, memory, attention, categorization, reasoning, causal inference, decision making, theory of mind
- Learning concepts from examples
- Learning causation from correlation
- Learning and applying intuitive theories (balancing complexity vs. fit)

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

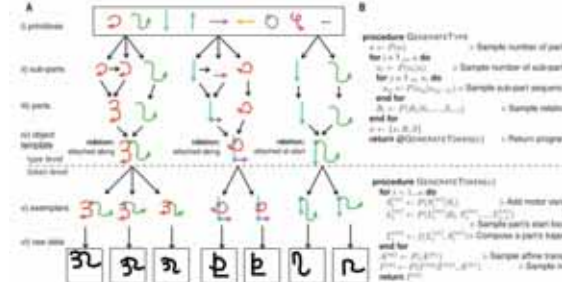


Tenenbaum, J. B., Griffiths, T. L. & Kemp, C. 2006. Theory-based Bayesian models of inductive learning and reasoning. Trends in cognitive sciences, 10, (7), 309-318.



Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. 2015. Human-level concept learning through probabilistic program induction. Science, 350, (6266), 1332-1338, doi:10.1126/science.aab3050.

A Bayesian program learning (BPL) framework, capable of learning a large class of visual concepts from just a single example and generalizing in ways that are mostly indistinguishable from people

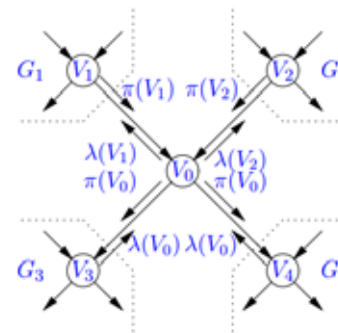


Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. 2015. Human-level concept learning through probabilistic program induction. Science, 350, (6266), 1332-1338, doi:10.1126/science.aab3050.

4) Graphs=Networks



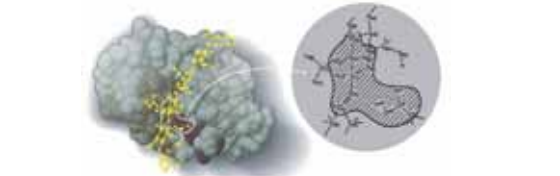
Image from <https://people.kth.se/~carlofi/teaching/FEL3250-2013/courseinfo.html>



Pearl, J. 1988. Embracing causality in default reasoning. Artificial Intelligence, 35, (2), 259-271.

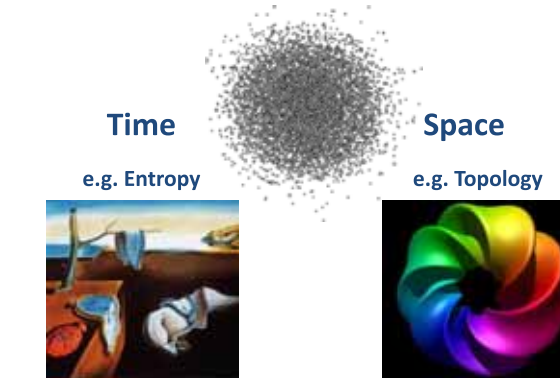


http://amturing.acm.org/vp/pearl_2658896.cfm



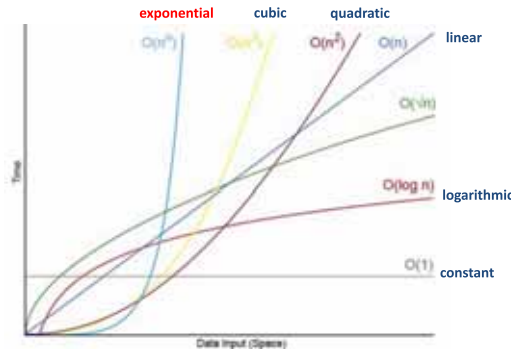
- Graphs as models for networks
- given as direct input (point cloud data sets)
- Given as properties of a structure
- Given as a representation of information (e.g. Facebook data, viral marketing, etc., ...)
- Graphs as nonparametric basis
- we learn the structure from samples and infer
- flat vector data, e.g. similarity graphs
- encoding structural properties (e.g. smoothness, independence, ...)

We skip this interesting chapter for now ...

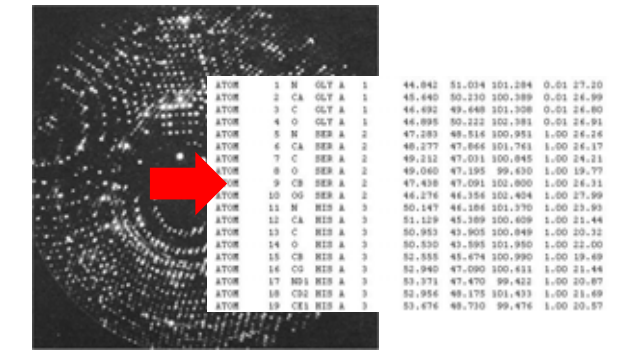


Dali, S. (1931) The persistence of memory

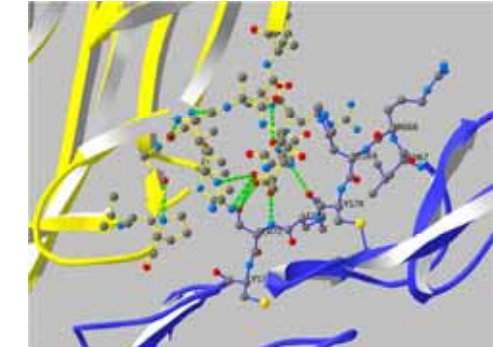
Bagula & Bourke (2012) Klein-Bottle



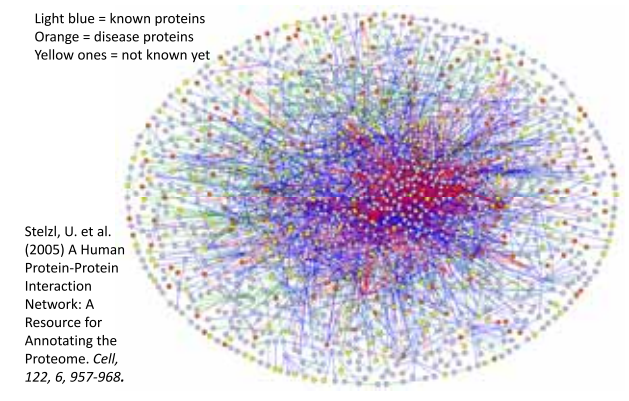
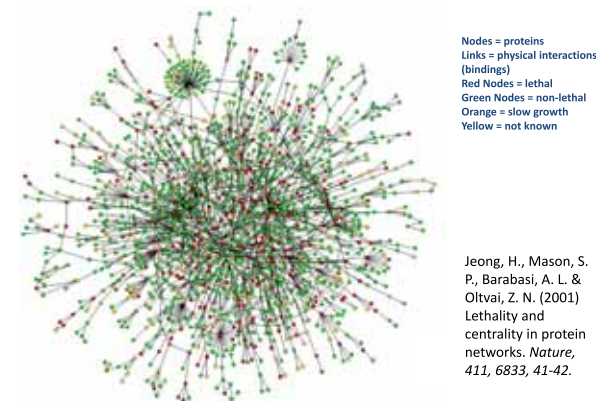
P versus NP and the Computational Complexity Zoo, please have a look at <https://www.youtube.com/watch?v=YX40hbAHx3s>



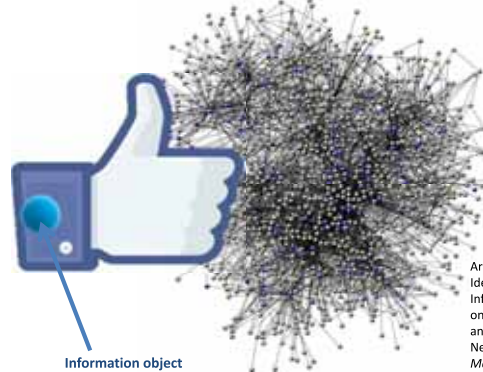
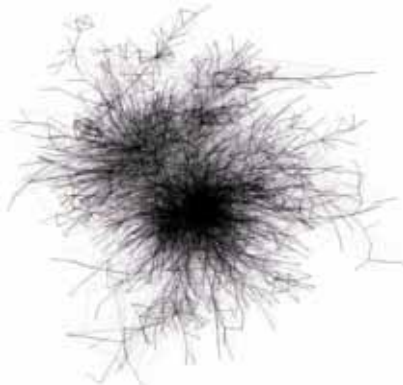
Wiltgen, M. & Holzinger, A. (2005) Visualization in Bioinformatics: Protein Structures with Physicochemical and Biological Annotations. In: *Central European Multimedia and Virtual Reality Conference, Prague, Czech Technical University (CTU), 69-74*



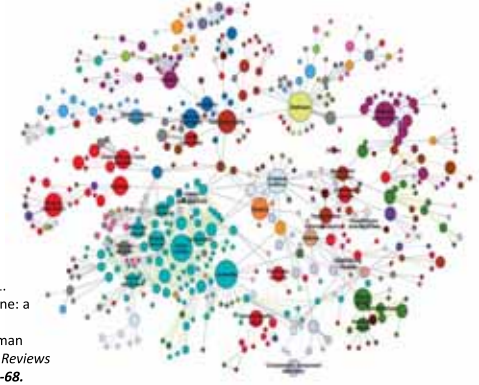
Wiltgen, M., Holzinger, A. & Tilz, G. P. (2007) Interactive Analysis and Visualization of Macromolecular Interfaces Between Proteins. In: *Lecture Notes in Computer Science (LNCS 4799)*. Berlin, Heidelberg, New York, Springer, 199-212.



Hurst, M. (2007), Data Mining: Text Mining, Visualization and Social Media. Online available: http://datamining.typepad.com/data_mining/2007/01/the_blogosphere.html, last access: 2011-09-24



Aral, S. (2011) Identifying Social Influence: A Comment on Opinion Leadership and Social Contagion in New Product Diffusion. *Marketing Science*, 30, 2, 217-223.



Barabási, A. L., Gulbahce, N. & Loscalzo, J. 2011. Network medicine: a network-based approach to human disease. *Nature Reviews Genetics*, 12, 56-68.

05 Bayesian Networks “Bayes’ Nets”

- is a **probabilistic model**, consisting of two parts:
- 1) a dependency structure and
- 2) local probability models.

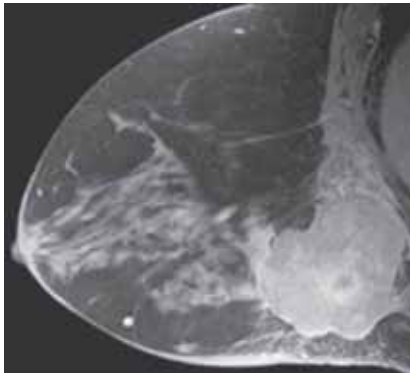
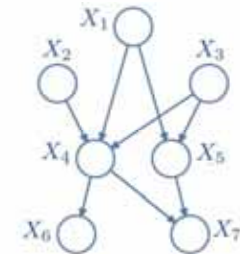
$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | Pa(x_i))$$

Where $Pa(x_i)$ are the parents of x_i

BN inherently model the uncertainty in the data. They are a successful marriage between probability theory and graph theory; allow to model a multidimensional probability distribution in a sparse way by searching independency relations in the data. Furthermore this model allows different strategies to integrate two data sources.

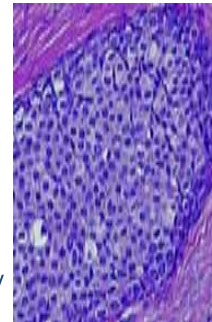
Pearl, J. (1988) *Probabilistic reasoning in intelligent systems: networks of plausible inference*. San Francisco, Morgan Kaufmann.

$$p(X_1, \dots, X_7) = p(X_1)p(X_2)p(X_3)p(X_4|X_1, X_2, X_3) \cdot p(X_5|X_1, X_3)p(X_6|X_4)p(X_7|X_4, X_5)$$

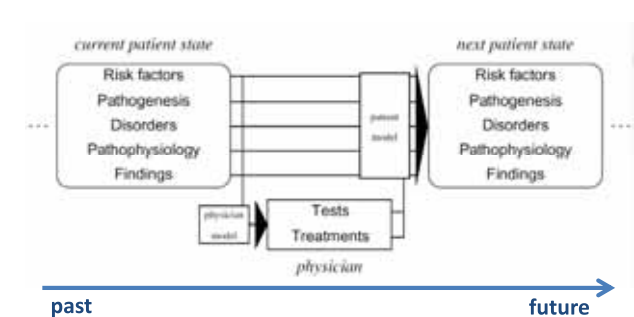


Overmoyer, B. A., Lee, J. M. & Lerwill, M. F. (2011) Case 17-2011 A 49-Year-Old Woman with a Mass in the Breast and Overlying Skin Changes. *New England Journal of Medicine*, 364, 23, 2246-2254.

- = the prediction of the future course of a disease conditional on the patient's history and a projected treatment strategy
- Danger: probable Information !
- Therefore valid prognostic models can be of great benefit for clinical decision making and of great value to the patient, e.g., for notification and quality of-life decisions



Knaus, W. A., Wagner, D. P. & Lynn, J. (1991) Short-term mortality predictions for critically ill hospitalized adults: science and ethics. *Science*, 254, 5030, 389.



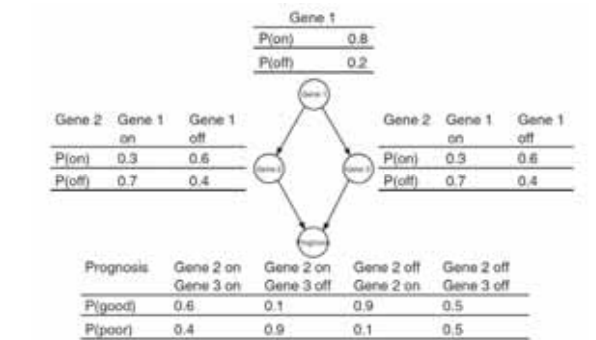
past

future

van Gerven, M. A. J., Taal, B. G. & Lucas, P. J. F. (2008) Dynamic Bayesian networks as prognostic models for clinical patient management. *Journal of Biomedical Informatics*, 41, 4, 515-529.

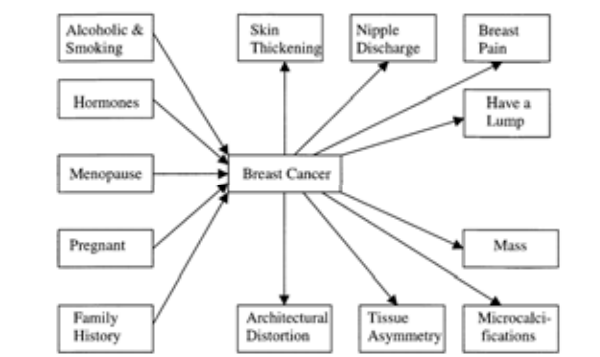
Category	Node description	State description
Diagnosis	Breast cancer	Present, absent.
Clinical history	Habit of drinking alcoholic beverages and smoking	Yes, no
	Taking female hormones	Yes, no
	Have gone through menopause	Yes, no
	Have ever been pregnant	Yes, no
	Family member has breast cancer	Yes, no
Physical findings	Nipple discharge	Yes, no
	Skin thickening	Yes, no
	Breast pain	Yes, no
	Have a lump(s)	Yes, no
Mammographic findings	Architectural distortion	Present, absent.
	Mass	Score from one to three, score from four to five, absent
	Microcalcification cluster	Score from one to three, score from four to five, absent
	Asymmetry	Present, absent.

Wang, X. H., et al. (1999) Computer-assisted diagnosis of breast cancer using a data-driven Bayesian belief network. *International Journal of Medical Informatics*, 54, 2, 115-126.

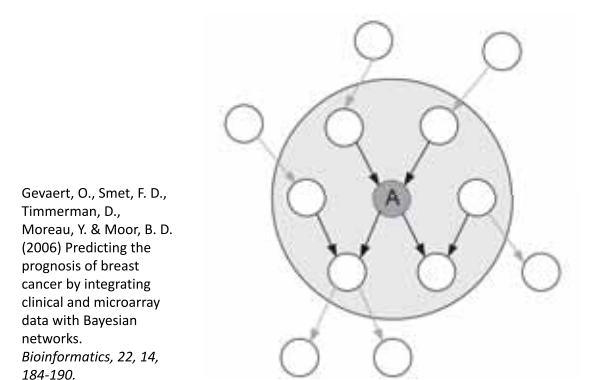


Gevaert, O., Smet, F. D., Timmerman, D., Moreau, Y. & Moor, B. D. (2006) Predicting the prognosis of breast cancer by integrating clinical and microarray data with Bayesian networks. *Bioinformatics*, 22, 14, 184-190.

- Next, N_{ij} is calculated by summing over all states of a variable:
- $N_{ij} = \sum_{k=1}^{r_i} N_{ijk} \cdot N'_{ijk}$ and N'_{ij} have similar meanings but refer to prior knowledge for the parameters.
- When no knowledge is available they are estimated using $N_{ijk} = N / (r_i q_i)$
- with N the equivalent sample size,
- r_i the number of states of variable i and
- q_i the number of instantiations of the parents of variable i .
- $\Gamma(\cdot)$ corresponds to the gamma distribution.
- Finally $p(S)$ is the prior probability of the structure.
- $p(S)$ is calculated by:
- $p(S) = \prod_{i=1}^n \prod_{l_i=1}^{p_i} p(l_i \rightarrow x_i) \prod_{m_i=1}^{o_i} p(m_i x_i)$
- with p_i the number of parents of variable x_i and o_i all the variables that are not a parent of x_i .
- Next, $p(a \rightarrow b)$ is the probability that there is an edge from a to b while $p(ab)$ is the inverse, i.e. the probability that there is no edge from a to b



Wang, X. H., et al. (1999) Computer-assisted diagnosis of breast cancer using a data-driven Bayesian belief network. *International Journal of Medical Informatics*, 54, 2, 115-126.



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.

- Estimating the parameters of the local probability models corresponding with the dependency structure.
- CPTs are used to model these local probability models.
- For each variable and instantiation of its parents there exists a CPT that consists of a set of parameters.
- Each set of parameters was given a uniform Dirichlet prior:

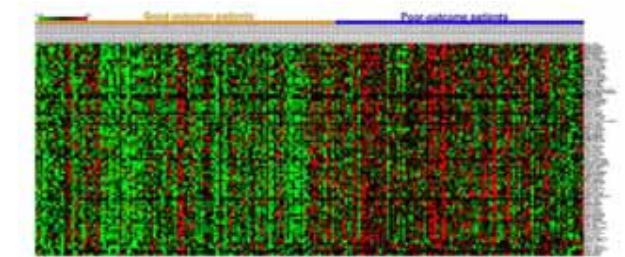
$$p(\theta_{ij}|S) = \text{Dir}(\theta_{ij}|N'_{ij1}, \dots, N'_{ijk}, \dots, N'_{ijr_i})$$

Note: With θ_{ij} a parameter set where i refers to the variable and j to the j -th instantiation of the parents in the current structure. θ_{ij} contains a probability for every value of the variable x_i given the current instantiation of the parents. Dir corresponds to the Dirichlet distribution with $(N'_{ij1}, \dots, N'_{ijr_i})$ as parameters of this Dirichlet distribution. Parameter learning then consists of updating these Dirichlet priors with data. This is straightforward because the multinomial distribution that is used to model the data, and the Dirichlet distribution that models the prior, are conjugate distributions. This results in a Dirichlet posterior over the parameter set:

$$p(\theta_{ij}|D, S) = \text{Dir}(\theta_{ij}|N'_{ij1} + N_{ij1}, \dots, N'_{ijk} + N_{ijk}, \dots, N'_{ijr_i} + N_{ijr_i})$$

with N_{ijk} defined as before.

- Integrating microarray data from multiple studies to increase sample size;
- = approach to the development of more robust prognostic tests

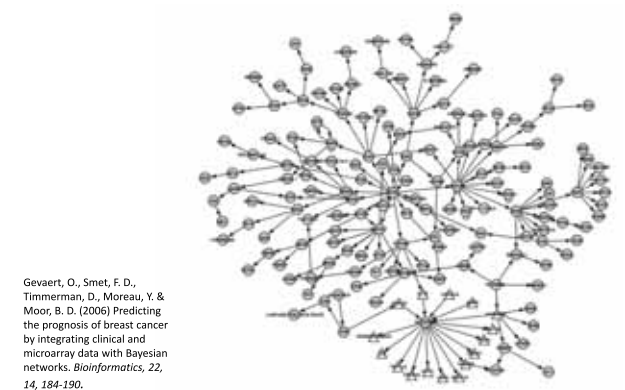


Xu, L., Tan, A., Winslow, R. & Geman, D. (2008) Merging microarray data from separate breast cancer studies provides a robust prognostic test. *BMC Bioinformatics*, 9, 1, 125-139.

- First the structure is learned using a [search strategy](#).
- Since the number of possible structures increases super exponentially with the number of variables,
- the well-known [greedy search algorithm K2](#) can be used in combination with the [Bayesian Dirichlet \(BD\) scoring metric](#):

$$p(S|D) \propto p(S) \prod_{i=1}^n \prod_{j=1}^{q_i} \left[\frac{\Gamma(N'_{ij})}{\Gamma(N'_{ij} + N_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N'_{ijk} + N_{ijk})}{\Gamma(N'_{ijk})} \right]$$

N_{ijk} ... number of cases in the data set D having variable i in state k associated with the j -th instantiation of its parents in current structure S .
 n is the total number of variables.



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.

- 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

Often it is better to have a good solution within time – than an perfect solution (much) later ...



06 Probabilistic Programming

- C → Probabilistic-C
- Scala → Figaro
- Scheme → Church
- Excel → Tabular
- Prolog → Problog
- Javascript → webPP
- Venture
- Python → PyMC



Probabilistic Program	Graphical Model
Variables	Variable nodes
Functions/operators	Factor nodes/edges
Fixed size loops/arrays	Plates
If statements	Gates (Minka & Winn)
Variable sized loops, Complex indexing, jagged arrays, mutation, recursion, objects/properties...	No common equivalent

1. Simple example: Nucleotide "A" may follow nucleotide "T" in the sequence more frequently for outcome X than for outcome Y.

2. $P(A|T, X) > P(A|T, Y)$

3. Specify the prior distribution: $P(\theta) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$

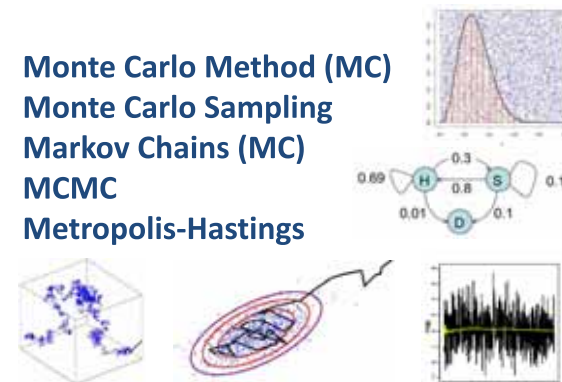
4. Specify the experimental data: $P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$

5. $P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$

6. Image Source: Dan Williams, Life Technologies, Austin TX

07 Markov Chain Monte Carlo (MCMC)

Monte Carlo Method (MC)
Monte Carlo Sampling
Markov Chains (MC)
MCMC
Metropolis-Hastings



- often we want to calculate characteristics of a **high-dimensional** probability distribution ... $p(\mathcal{D}|\theta)$

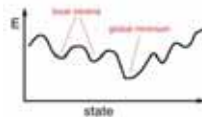
$$p(h|d) \propto p(\mathcal{D}|\theta) * p(h)$$

Posterior integration problem: (almost) all statistical inference can be deduced from the posterior distribution by calculating the appropriate sums, which involves an integration:

$$J = \int f(\theta) * p(\theta|\mathcal{D}) d\theta$$



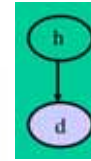
- Solving intractable integrals
- Bayesian statistics: **normalizing** constants, expectations, marginalization
- Stochastic Optimization
- Generalization of simulated annealing
- Monte Carlo expectation maximization (EM)



- Statistical physics:** computing the partition function – this is evaluating the posterior probability of a hypothesis and this requires summing over all hypotheses ... remember:

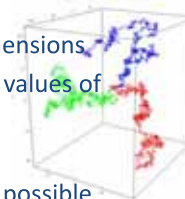
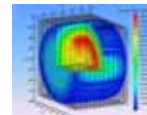
$$\mathcal{H} = \{H_1, H_2, \dots, H_n\} \quad \forall (h, d)$$

$$P(h|d) = \frac{P(d|h) * P(h)}{\sum_{h' \in \mathcal{H}} P(d|h') P(h')}$$



- Class of algorithms that rely on **repeated random sampling**
- Basic idea: using **randomness** to solve problems with high uncertainty (Laplace, 1781)
- For solving **multidimensional integrals** which would otherwise intractable
- For simulation of systems with **many dof**
- e.g. fluids, gases, particle collectives, **cellular structures** - see our last tutorial on Tumor growth simulation!

- Physical simulation
- estimating neutron diffusion time
- Computing expected utilities and best responses toward Nash equilibria
- Computing volumes in high-dimensions
- Computing eigen-functions and values of operators (e.g. Schrödinger)
- Statistical physics
- Counting many things as fast as possible



- for solving problems of probabilistic inference involved in developing computational models
- as a source of hypotheses about how the human mind might solve problems of inference
- For a function $f(x)$ and distribution $P(x)$, the expectation of f with respect to P is generally the average of f , when x is drawn from the probability distribution $P(x)$

$$\mathbb{E}_{p(x)}(f(x)) = \sum_x f(x) P(x) dx$$

- Expectation of a function $f(x, y)$ with respect to a random variable x is denoted by $\mathbb{E}_x[f(x, y)]$
- In situations where there is no ambiguity as to which variable is being averaged over, this will be simplified by omitting the suffix, for instance $\mathbb{E}x$.
- If the distribution of x is conditioned on another variable z , then the corresponding conditional expectation will be written $\mathbb{E}_x[f(x)|z]$
- Similarly, the variance is denoted $var[f(x)]$, and for vector variables the covariance is written $cov[x, y]$

$$\operatorname{argmax}_x f(x)$$

Normalization: $p(x|y) = \frac{p(y|x) * p(x)}{\int_X p(y|x) * p(x) dx}$

Marginalization: $p(x) = \int_Z p(x, z) dz$

Expectation: $\mathbb{E}_{p(x)}(f(x)) = \int_X f(x)p(x)dx$

08 Metropolis-Hastings Algorithm

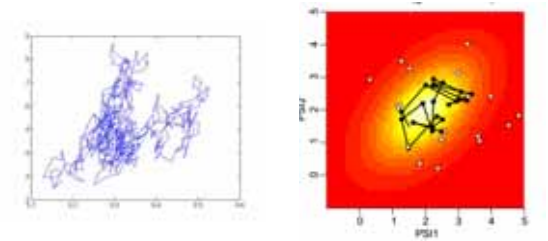


Image Source: Peter Mueller, Anderson Cancer Center

Barber, D. 2012. Bayesian reasoning and machine learning, Cambridge, Cambridge University Press, p. 500

```

1: Choose a starting point  $x^1$ .
2: for  $i = 2$  to  $L$  do
3:   Draw a candidate sample  $x^{cand}$  from the proposal  $\tilde{q}(x^i|x^{i-1})$ .
4:   Let  $\alpha = \frac{\tilde{q}(x^{i-1}|x^{cand})p(x^{cand})}{\tilde{q}(x^{cand}|x^{i-1})p(x^{i-1})}$ .
5:   if  $\alpha \geq 1$  then  $x^i = x^{cand}$ 
6:   else
7:     draw a random value  $u$  uniformly from the unit interval  $[0, 1]$ .
8:     if  $u < \alpha$  then  $x^i = x^{cand}$ 
9:     else  $x^i = x^{i-1}$ 
10:  end if
11: end if
12: end if
13: end for
    
```



- Importance sampling is a technique to approximate averages with respect to an intractable distribution $p(x)$.
- The term 'sampling' is arguably a misnomer since the method does not attempt to draw samples from $p(x)$.
- Rather the method draws samples from a simpler importance distribution $q(x)$ and then reweights them
- such that averages with respect to $p(x)$ can be approximated using the samples from $q(x)$.

- The Gibbs Sampler is an interesting special case of MH:

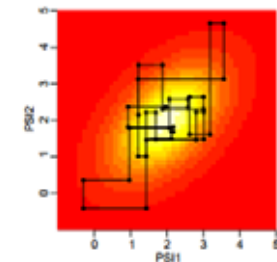
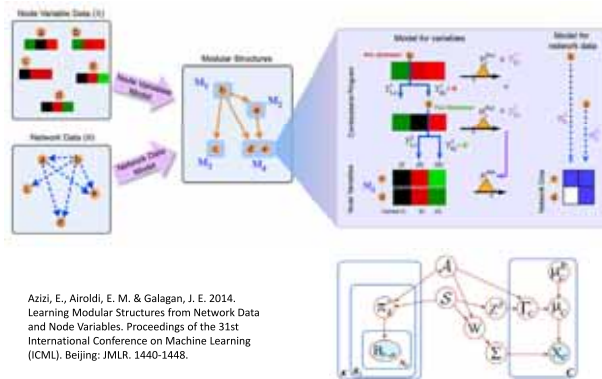


Image Source: Peter Mueller, Anderson Cancer Center



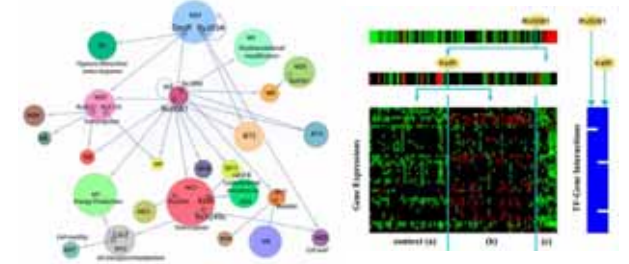
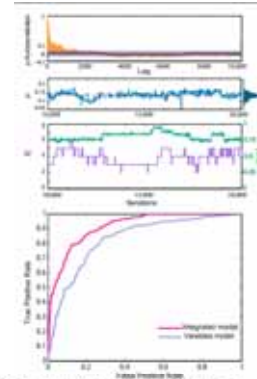
Azizi, E., Airolidi, E. M. & Galagan, J. E. 2014. Learning Modular Structures from Network Data and Node Variables. Proceedings of the 31st International Conference on Machine Learning (ICML). Beijing: JMLR. 1440-1448.

Algorithm 1 RIMCMC for sampling parameters

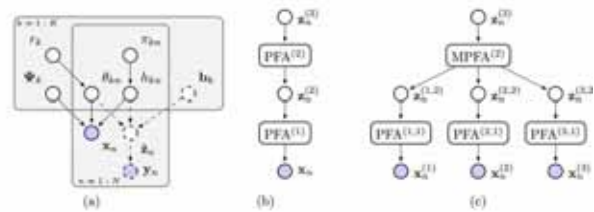
Inputs:
Node Variables Data X ,
Network Data B

for iterations $j = 1$ to J do
 Sample $\mathcal{A}^{(j+1)}$ given $\mathcal{A}^{(j)}$ using Alg 2 in (Azizi et al., 2014).
 Sample $\Sigma^{(j+1)}$ given $\mathcal{A}^{(j+1)}$ using Alg 3 in (Azizi et al., 2014).
 for modules $k = 1$ to $K^{(j+1)}$ do
 Propose $\alpha_k^{(j+1)} \sim \mathcal{N}(\alpha_k^{(j)}, F)$.
 Accept with probability P_{acc} ; update $\Sigma^{(j+1)}$.
 end for
 for parents $r = 1$ to R_k do
 Propose $\alpha_r^{(j+1)} \sim \mathcal{N}(\alpha_r^{(j)}, F)$; accept with P_{acc} .
 Propose $\alpha_r^{(j+1)} \sim \mathcal{N}(\alpha_r^{(j)}, F)$; accept with P_{acc} .
 end for
 for condition $c = 1$ to C do
 Propose $\mu_c^{(j+1)} \sim \mathcal{N}(\mu_c^{(j)}, F)$; accept with P_{acc} .
 Propose $\gamma_c^{(j+1)} \sim \mathcal{N}(\gamma_c^{(j)}, F)$; accept with P_{acc} .
 end for
end for

Azizi, E., Airolidi, E. M. & Galagan, J. E. 2014. Learning Modular Structures from Network Data and Node Variables. Proceedings of the 31st International Conference on Machine Learning (ICML). Beijing: JMLR. 1440-1448.



Azizi, E., Airolidi, E. M. & Galagan, J. E. 2014. Learning Modular Structures from Network Data and Node Variables. Proceedings of the 31st International Conference on Machine Learning (ICML). Beijing: JMLR. 1440-1448.



Henao, R., Lu, J. T., Lucas, J. E., Ferranti, J. & Carin, L. 2016. Electronic health record analysis via deep poisson factor models. *Journal of Machine Learning Research JMLR*, 17, 1-32.



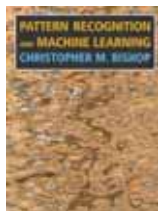
Thank you!

Questions

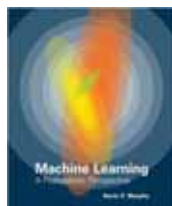
- What is the main difference between the ideas of Pierre Simon de Laplace and Lady Lovelace?
- What is medical action consisting most of the time?
- How does a human make a decision - as far as we know?
- What is the main idea of a probabilistic programming language?
- Why did Judea Pearl receive the Turing Award (Noble Prize in Computer Science)?
- What fields are coming together in PGM?
- What are the challenges in network structures?
- Give a classification of Graphical Models!
- What are plates and nested plates?
- Provide corresponding examples of metabolic networks!

- What is a factored graph?
- Describe the protein structure prediction problem! Why is it hard?
- Why are protein-protein interactions so important?
- Describe the problem of graph-isomorphism!
- How does a Bayes Net work?
- Why is predicting important in clinical medicine?
- What is a Markov-Blankett?
- Which two tasks do we have in Graphical Model Learning?
- Why would we need probabilistic programming languages?
- Describe the main idea of MCMC!
- What is the main problem in marginalization?
- What is the benefit of the MH Algorithm?

Appendix



Bishop, C. M. 2007. *Pattern Recognition and Machine Learning*, Heidelberg, Springer. Chapter 8 on graphical models openly available: <http://research.microsoft.com/en-us/um/people/cmbishop/prml/>



Murphy, K. P. 2012. *Machine learning: a probabilistic perspective*, MIT press. Chapter 26 (pp. 907) – Graphical model structure learning



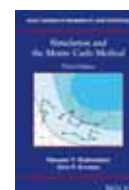
Koller, D. & Friedman, N. 2009. *Probabilistic graphical models: principles and techniques*, MIT press.



Rubinstein, R. Y. & Kroese, D. P. 2013. *The cross-entropy method: a unified approach to combinatorial optimization, Monte-Carlo simulation and machine learning*, Springer



Cameron Davidson-Pilon 2015. *Bayesian methods for hackers: probabilistic programming and Bayesian inference*, Addison-Wesley Professional.



Rubinstein, R. Y. & Kroese, D. P. 2013. *Simulation and the Monte-Carlo Method*, Wiley



Stiller, A., Goodman, N. & Frank, M. C. Ad-hoc scalar implicature in adults and children. *CogSci*, 2011.



<https://goo.gl/6a7rOC>

Chapter 8 Graphical Models is as sample chapter fully downloadable for free

Bishop, C. M. 2006. Pattern Recognition and Machine Learning, Heidelberg, Springer.



<http://bayes.cs.ucla.edu/BOOK-2K/>

Pearl, J. 2009. Causality: Models, Reasoning, and Inference (2nd Edition), Cambridge, Cambridge University Press.

Remember: Three main types of Machine Learning

- I) Supervised learning (classification)
 - $y = f(x)$
 - Given x , y pairs; find a f that map a new x to a proper y
 - Regression, logistic regression, classification
 - Expert provides examples e.g. classification of clinical images
 - Disadvantage: Supervision can be expensive
- II) Unsupervised learning (clustering)
 - $f(x)$
 - Given x (features only), find f that gives you a description of x
 - Find similar points in high-dim X
 - E.g. clustering of medical images based on their content
 - Disadvantage: Not necessarily task relevant
- III) Reinforcement learning
 - $y = f(x)$
 - more general than supervised/unsupervised learning
 - learn from interaction to achieve a goal
 - Learning by direct interaction with environment (automatic ML)
 - Disadvantage: broad difficult approach, problem with high-dim data

12,081 as of 10.4.2018 - 10,624 citations 26.03.2017

Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. Biometrika, 57, (1), 97-109.



- Medicine is an extremely complex application domain – dealing most of the time with uncertainties -> **probable information!**
- Key: Structure learning and prediction in large-scale biomedical networks with probabilistic graphical models
- Causality and Probabilistic Inference
- Uncertainties are present at all levels in health related systems
- Data sets from which ML learns are noisy, mislabeled, atypical, etc. etc.
- Even with data of high quality, gauging and combining a multitude of data sources and constraints in usually imperfect models of the world requires us to represent and process **uncertain knowledge** in order to make **viable decisions in context and within reasonable time!**
- In the increasingly complicated settings of modern science, model structure or causal relationships may not be known a-priori [1].
- Approximating probabilistic inference in Bayesian belief networks is NP-hard [2] -> here we need the “human-in-the-loop” [3]

[1] Sun, X., Janzing, D. & Schölkopf, B. Causal Inference by Choosing Graphs with Most Plausible Markov Kernels. IJCAI, 2006.

[2] Dagum, P. & Luby, M. 1993. Approximating probabilistic inference in Bayesian belief networks is NP-hard. Artificial Intelligence, 60, (1), 141-153.

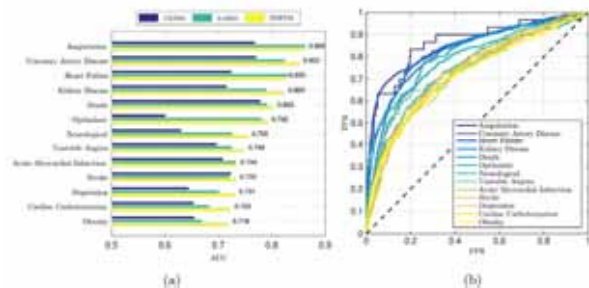
[3] Holzinger, A. 2016. Interactive Machine Learning for Health Informatics: When do we need the human-in-the-loop? Springer Brain Informatics (BRIN), 3, 1-13. doi:10.1007/s40708-016-0042-6.

5,223 citations as of 26.03.2017



Image Source:
<http://www.manhattanprojectvoices.org/oral-histories/nicholas-metropolis-interview>

MCMC based DPFM outperforms other approaches



Henao, R., Lu, J. T., Lucas, J. E., Ferranti, J. & Carin, L. 2016. Electronic health record analysis via deep poisson factor models. Journal of Machine Learning Research JMLR, 17, 1-32.

- Reinforcement Learning is the **oldest approach**, with the longest history and can provide insight into understanding human learning [1]
- RL is the **“AI problem in the microcosm”** [2]
- Future opportunities are in Multi-Agent RL (MARL), Multi-Task Learning (MTL), Generalization and Transfer-Learning [3], [4].

[1] Turing, A. M. 1950. Computing machinery and intelligence. Mind, 59, (236), 433-460.

[2] Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. Nature, 521, (7553), 445-451. doi:10.1038/nature14540.

[3] Taylor, M. E. & Stone, P. 2009. Transfer learning for reinforcement learning domains: A survey. The Journal of Machine Learning Research, 10, 1633-1685.

[4] Pan, S. J. & Yang, Q. A. 2010. A Survey on Transfer Learning. IEEE Transactions on Knowledge and Data Engineering, 22, (10), 1345-1359. doi:10.1109/tkde.2009.191.

37,202 (as of 10.4.2018) - 34,140 citations (26.03.2017)



Finally a practical example

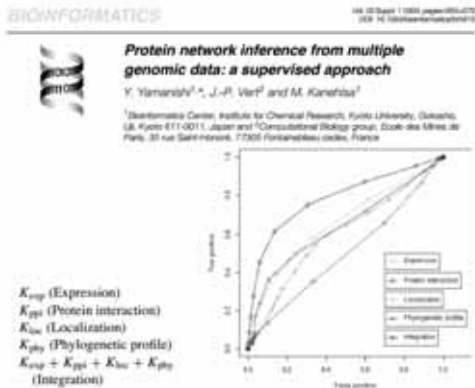
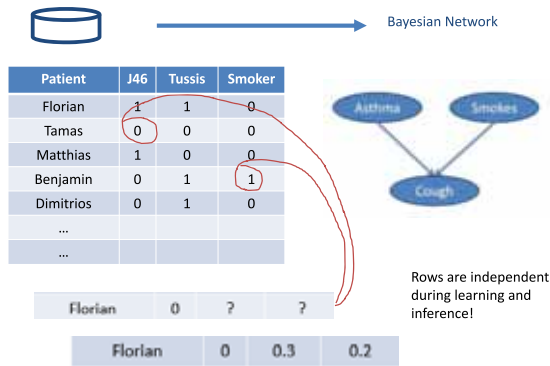
Graphical Model Learning

- Remember: GM are a marriage between probability theory and graph theory and provide a tool for dealing with our two grand challenges in the biomedical domain:

Uncertainty and complexity

- The learning task is two-fold:
 - 1) Learning unknown probabilities
 - 2) Learning unknown structures

Jordan, M. I. 1998. Learning in graphical models, Springer



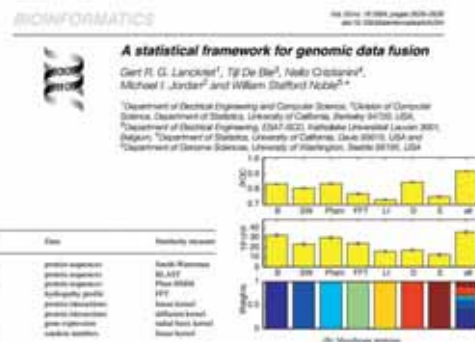
- 1) Test if a distribution is decomposable with regard to a given graph.
 - This is the most direct approach. It is not bound to a graphical representation,
 - It can be carried out w.r.t. other representations of the set of subspaces to be used to compute the (candidate) decomposition of a given distribution.
- 2) Find a suitable graph by measuring the strength of dependences.
 - This is a heuristic, but often highly successful approach, which is based on the frequently valid assumption that in a conditional independence graph an attribute is more strongly dependent on adjacent attributes than on attributes that are not directly connected to them.
- 3) Find an independence map by conditional independence tests.
 - This approach exploits the theorems that connect conditional independence graphs and graphs that represent decompositions.
 - It has the advantage that a single conditional independence test, if it fails, can exclude several candidate graphs. Beware, because wrong test results can thus have severe consequences.

Borgelt, C., Steinbrecher, M. & Kruse, R. R. 2009. Graphical models: representations for learning, reasoning and data mining, John Wiley & Sons.

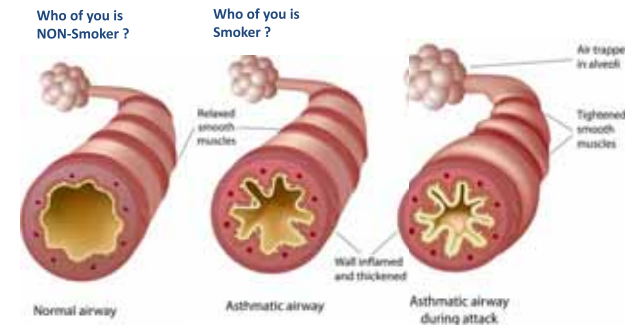
- Asthma can be hereditary
- Friends may have similar smoking habits
- Augmenting graphical model with relations between the entities – Markov Logic



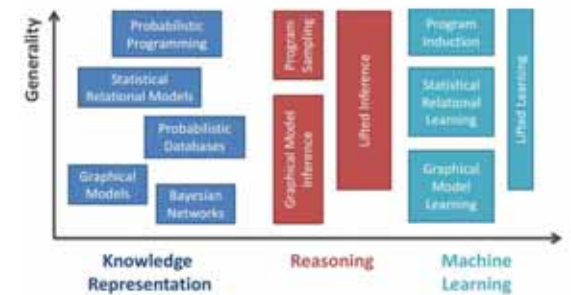
- 2.1 Asthma \Rightarrow Cough
- 3.5 Smokes \Rightarrow Cough
- 2.1 Asthma(x) \Rightarrow Cough(x)
- 3.5 Smokes(x) \Rightarrow Cough(x)
- 1.9 Smokes(x) \wedge Friends(x,y) \Rightarrow Smokes(y)
- 1.5 Asthma(x) \wedge Family(x,y) \Rightarrow Asthma(y)



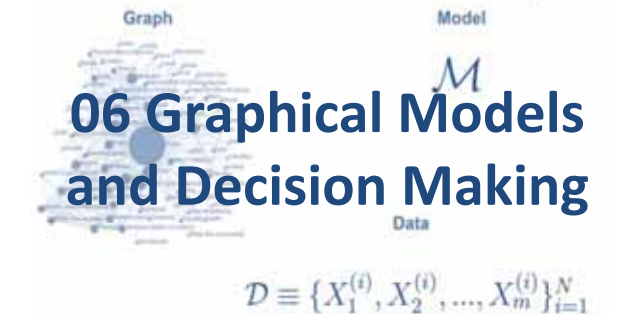
Lanckriet, G. R., De Bie, T., Cristianini, N., Jordan, M. I. & Noble, W. S. 2004. A statistical framework for genomic data fusion. Bioinformatics, 20, (16), 2626-2635.

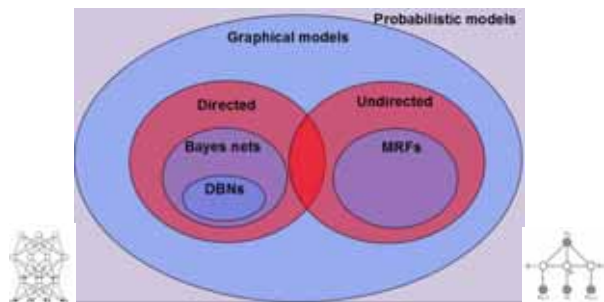


Beasley, R. 1998. Worldwide variation in prevalence of symptoms of asthma, allergic rhinoconjunctivitis, and atopic eczema: ISAAC. The Lancet, 351, (9111), 1225-1232, doi:http://dx.doi.org/10.1016/S0140-6736(97)07302-9.



Example for probabilistic rule learning, in which probabilistic rules are learned from probabilistic examples: The ProbFOIL+ Algorithm solves this problem by combining the principles of the rule learner FOIL with the probabilistic Prolog called Problog, see: De Raedt, L., Dries, A., Thon, I., Van Den Broeck, G. & Verbeke, M. 2015. Inducing probabilistic relational rules from probabilistic examples. International Joint Conference on Artificial Intelligence (IJCAI).



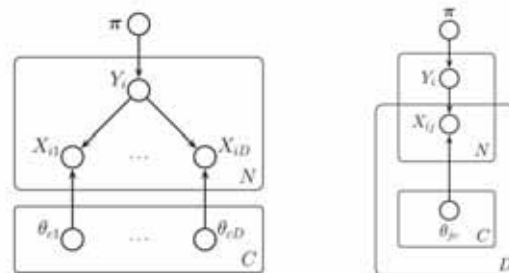


Murphy, K. P. 2012. Machine learning: a probabilistic perspective, Cambridge (MA), MIT press.

Holzinger Group hci-kdd.org

127

MAKE Health Module 04



$\pi \dots$ multinomial parameter vector, Stationary distribution of Markov chain

Holzinger Group hci-kdd.org

128

MAKE Health Module 04

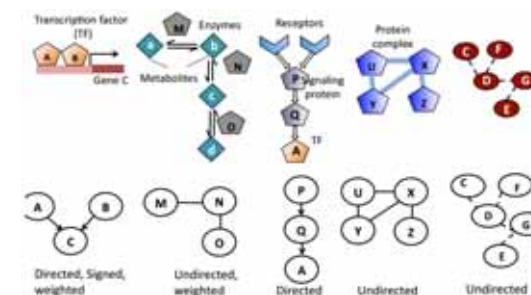


Image credit to Anna Goldenberg, Toronto

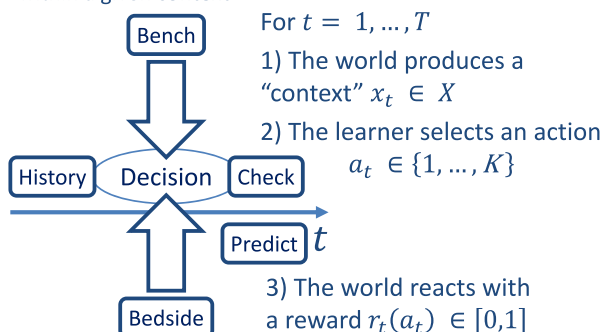
Holzinger Group hci-kdd.org

129

MAKE Health Module 04

TU Decision Making: Learn good policy for selecting actions

Goal: Learn an **optimal policy** for selecting best actions within a given **context**



Holzinger Group hci-kdd.org

130

MAKE Health Module 04

TU GM are amongst the most important ML developments

Key Idea: Conditional independence assumptions are very useful – however: Naïve Bayes is extreme!

X is *conditionally independent* of Y , given Z , if the $P(X)$ governing X is independent of value Y , given value of Z :

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

can be abbr. with $P(X|Y, Z) = P(X|Z)$

- Graphical models express sets of conditional independence assumptions via graph structure
- The graph structure plus associated parameters define joint probability distribution over the set of variables

Holzinger Group hci-kdd.org

131

MAKE Health Module 04

TU Remember

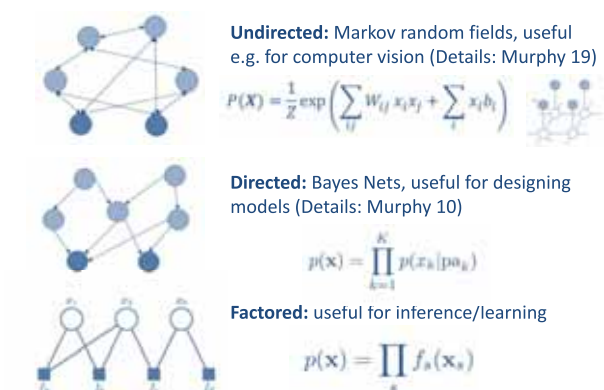
- Medicine is an extremely complex application domain – dealing most of the time with uncertainties -> **probable information!**
- When we have big data but little knowledge automatic ML can help to gain insight:
- Structure learning and prediction in large-scale biomedical networks with probabilistic graphical models**
- If we have little data and deal with NP-hard problems we still need the human-in-the-loop

Holzinger Group hci-kdd.org

132

MAKE Health Module 04

TU Three types of Probabilistic Graphical Models

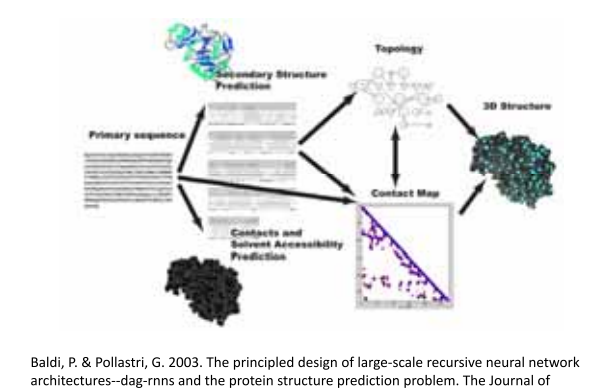


Holzinger Group hci-kdd.org

133

MAKE Health Module 04

TU From structure to function prediction



Holzinger Group hci-kdd.org

134

MAKE Health Module 04

TU Protein Network Inference

- Hypothesis: most biological functions involve the interactions between many proteins, and the complexity of living systems arises as a result of such interactions.
- In this context, the problem of inferring a global protein network for a given organism,
- using all (genomic) data of the organism,
- is one of the main challenges in computational biology

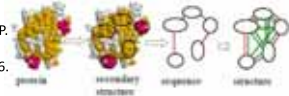
Yamanishi, Y., Vert, J.-P. & Kanehisa, M. 2004. Protein network inference from multiple genomic data: a supervised approach. Bioinformatics, 20, (suppl 1), i363-i370.

Holzinger Group hci-kdd.org

135

MAKE Health Module 04

Borgwardt, K. M., Ong, C. S., Schönaier, S., Vishwanathan, S., Smola, A. J. & Kriegel, H.-P. 2005. Protein function prediction via graph kernels. *Bioinformatics*, 21, (suppl 1), i47-i56.



- Important for health informatics: Discovering relationships between biological components
- Unsolved problem in computer science:
- Can the graph isomorphism problem be solved in polynomial time?
 - So far, no polynomial time algorithm is known.
 - It is also not known if it is NP-complete
 - We know that subgraph-isomorphism is NP-complete