## Slide 1

Andreas Holzinger
185.A83 Machine Learning for Health Informatics
2018S, VU, 2.0 h, 3.0 ECTS
Lecture 06 - Module 04 – Week 20 - 15.05.2018

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

a.holzinger@hci-kdd.org
http://hci-kdd.org/machine-learning-for-health-informatics-course
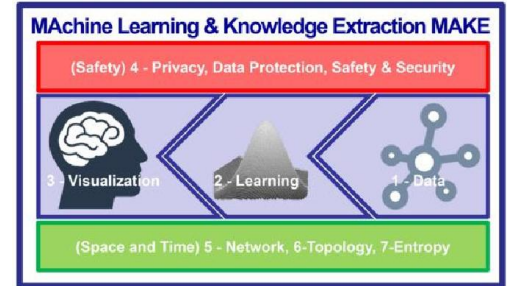
## Slide 2

**Science is to test crazy ideas –**
**Engineering is to put these ideas into Business**
**Lucky Students** ☺

## Slide 3

MAchine Learning & Knowledge Extraction MAKE
(Safety) 4 - Privacy, Data Protection, Safety & Security
3 - Visualization    2 - Learning    1 - Data
(Space and Time) 5 - Network, 6-Topology, 7-Entropy

## Slide 4

Cognition | Visualization | Data fusion
Perception | | Preprocessing
Decision | Interaction | Integration

CONCEPTS | THEORIES | PARADIGMS | MODELS | METHODS | TOOLS
Dimensionality | Complexity | Unsupervised | Gaussian P. | Regularization | Python
Reinforcement | Bayesian p(x) | Supervised | Graphical M. | Scaling | Church
Representation | Entropy/KL | Semi-Superv. | Neural Nets | Aggregation | Anglican
No-free-lunch | Vapnik-Chernov. | iML | Kernel/SVM | Evolution | Julia
Multi-Task Learning | Transfer Learning | Multi-Agent-Hybrid-Systems

Data Protection, Safety and Security and Privacy Aware Machine Learning (PAML)
Application, Validation, Evaluation, Impact – Social, Economic, Acceptance, Trust

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

## Slide 5

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

## Slide 6



00 Reflection

http://smashinghub.com/beautiful-examples-of-shadow-photography

## Slide 7

04

- 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.
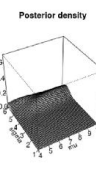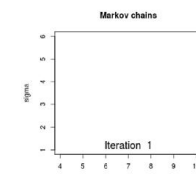- 6) **understand** the data in the **context** of an application domain

## Slide 8

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

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

Compute $a_i := \sum_j J_{ij}x_j$
Draw $u$ from Uniform(0,1)
If $u < 1/(1 + e^{-2a_i})$
  $x_i := +1$
Else
  $x_i := -1$
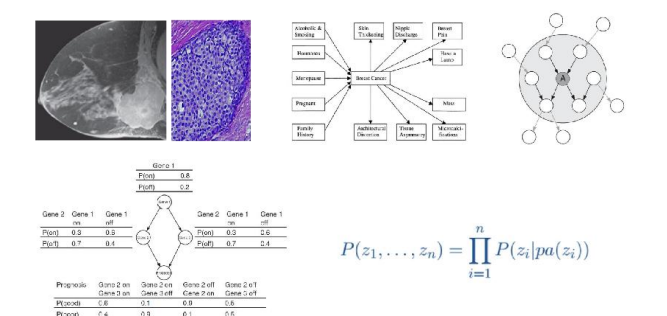
Markov chains          Posterior density

Propp, J. G. & Wilson, D. B. 1996. Exact sampling with coupled Markov chains and applications to statistical mechanics. Random structures and Algorithms, 9, (1-2), 223-252.

## Slide 9

$$P(z_1, \ldots, z_n) = \prod_{i=1}^{n} P(z_i | pa(z_i))$$

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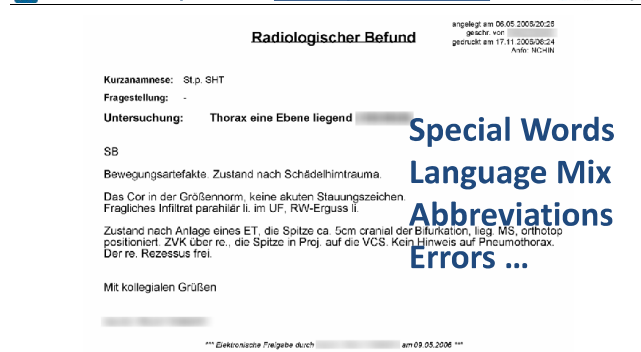
- For certain cases it is tractable if:
  - Just one variable is unobserved
  - We have singly connected graphs (no undirected loops -> belief propagation)
  - Assigning probability to fully observed set of variables
- Possibility: Monte Carlo Methods (generate many samples according to the Bayes Net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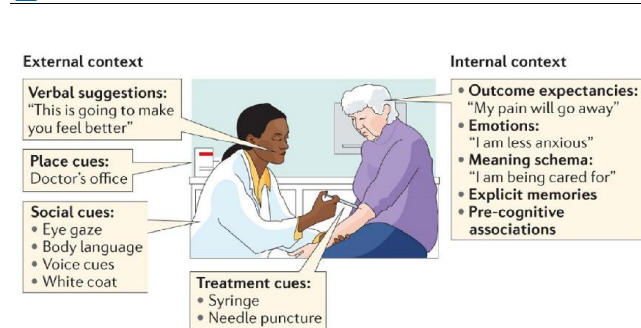- Otherwise: approximate solutions, NOTE:
  **Sometimes it is better to have an approximate solution to a complex problem – than a perfect solution to a simplified problem**

---

**Radiologischer Befund**
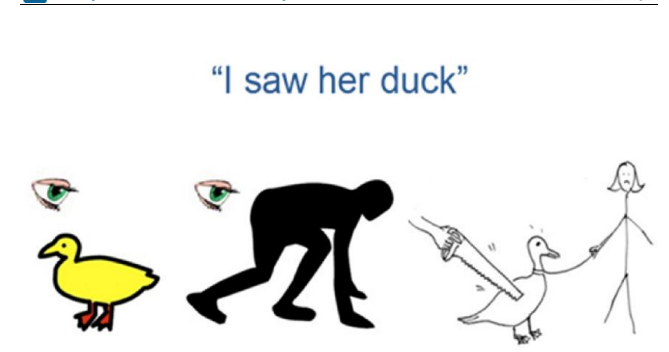
Special Words
Language Mix
Abbreviations
Errors …

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

---

"I saw her duck"

---

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

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

---

- Intelligence?
  - Hundreds of controversial definitions – very hard to define;
  - For us: ability to solve problems, to make decisions and to acquire and apply knowledge and skills
- Learning?
  - Different definitions – relatively hard to define
  - basically acquisition of knowledge through prior experience
- Problem Solving?
  - Process of finding solutions to complex issues
- Reasoning?
  - ability of our mind to think and understand things
- Sense Making?
  - Process of giving meaning to experience
- Causality?
  - Relationship between cause and effect
- Decision Making?
  - Process of "de-ciding" ("ent-scheiden") between alternative options

---

# 01 Probabilistic Decision Making

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

---

**Medical action …**

**is permanent decision making under uncertainty …**

Medical Decision Making
Harold C. Sox
Michael C. Higgins
Douglas K. Owens
SECOND EDITION
WILEY-BLACKWELL

---

Wager, T. D. & Atlas, L. Y. 2015. The neuroscience of placebo effects: connecting context, learning and health. Nat Rev Neurosci, 16, (7), 403-418, doi:10.1038/nrn3976

---

- **Type 1 Decisions:** related to the **diagnosis,** i.e. computers are used to assist in diagnosing a disease on the basis of the individual patient data. Questions include:
  - What is the probability that this patient has a myocardial infarction on the basis of given data (patient history, ECG, …)?
  - What is the probability that this patient has acute appendices, given the signs and symptoms concerning abdominal pain?
- **Type 2 Decisions:** related to **therapy,** i.e. computers are used to select the best therapy on the basis of clinical evidence, e.g.:
  - What is the best therapy for patients of age x and risks y, if an obstruction of z % is seen in the left coronary artery?
  - What amount of insulin should be prescribed for a patient during the next 5 days, given the blood sugar levels and the amount of insulin taken during the recent weeks?

Harold C. Sox, Michael C. Higgins & Douglas K. Owens 1988. Medical decision making, Second Edition, Chichester, Wiley.

## Slide 19

$$\mathbb{E}[f] = \int p(x)f(x)\, dx$$

$$\mathbb{E}[f] \simeq \frac{1}{N}\sum_{n=1}^{N} f(x_n)$$

## Slide 20

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

The expected utility of decision $D = d$ is

http://www.eoht.info/page/Oskar+Morgenstern

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

An optimal single decision is the decision $D = dmax$ whose expected utility is maximal:

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

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

## Slide 21

Bayesian Data Analysis
Third Edition

Andrew Gelman, John B. Carlin, David B. Dunson, Aki Vehtari and Donald B. Rubin

https://github.com/avehtari/BDA_py_demos

http://www.stat.columbia.edu/~gelman/book/data/

Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari & Donald B. Rubin 2014. Bayesian data analysis, Boca Raton (FL), CRC press.

## Slide 22

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

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

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



Recognize speech

Wreck a nice beach

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

## Slide 23

- 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)

## Slide 24

- 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.

## Slide 25

LTM: Prior knowledge $\mathcal{H}$

$$\mathcal{H} = \{H_1, H_2, \dots, H_n\}$$

STM:

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

Uncertain world

## Slide 26

```
for t = 1,…, n do
    The agent perceives state s_t
    The agent performs action a_t
    The environment evolves to s_{t+1}
    The agent receives reward r_t
end for
```

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



Agent
Representation
Learning algorithm
Action selection policy

State $x^{(t)}$

Reward $r^{(t)}$   $r^{(t+1)}$

Action $a^{(t)}$

Environment

$x^{(t+1)}$

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

## Slide 27

Belief   Desire   Action

?   a   b   c

{a,b,c}

→ decision that is best for worst case

Non-deterministic model

~ Adversarial search

?   a   b   c

{a($p_a$), b($p_b$), c($p_c$)}

→ **decision that maximizes expected utility value**

Probabilistic model

# 02 Probabilistic Programming

---

- Dan ROY: Probabilistic Programming Wiki
  http://www.probabilistic-programming.org/wiki/Home
- Frank WOOD, many tutorials, slides, code and papers
  http://www.robots.ox.ac.uk/~fwood/teaching/index.html
- Avi PFEFFER 2016. Practical probabilistic programming, Shelter Island (NY), Manning
  https://www.manning.com/books/practical-probabilistic-programming

  Look also for work of:
  Andrew GORDON
  Noah GOODMAN
  Josh TENENBAUM
  John WINN
  Rob ZINKOV
  Vikash MANSINGHA
  David WINGATE

---

- 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 …
- Prior = belief before making a particular observation
- Posterior = belief after making the observation and is the prior for the next observation – intrinsically incremental

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})} \qquad p(x_i|y_j) = \frac{p(y_j|x_i)p(x_i)}{\sum p(x_i, y_j)p(x_i)}$$

---

Observed data:

≈ Training data: $\mathcal{D} = x_{1:n} = \{x_1, x_2, ..., x_n\}$   $x, y$  $A, B, ...$

Feature Parameter: $\theta$ or hypothesis $h$   $h \in \mathcal{H}$

Prior belief ≈ prior probability of hypothesis $h$: $p(\theta)$   $p(h)$

Likelihood ≈ $p(x)$ of the data that $h$ is true   $p(\mathcal{D}|\theta)$   $p(d|h)$

Data evidence ≈ marginal $p(x)$ that $h$ = true   $p(\mathcal{D})$   $\sum_{h \in \mathcal{H}} p(d|h) * p(h)$

Posterior ≈ $p(x)$ of $h$ after seen ("learn") data $d$   $p(\theta|\mathcal{D})$   $p(h|d)$

$$posterior = \frac{likelihood * prior}{evidence} \qquad p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) * p(\theta)}{p(\mathcal{D})}$$

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

---

**Inference**

Parameters → Program → Output

$\theta \, p(\mathbf{x}|\mathbf{y})$ → $p(X|\theta)p(\theta)$ $p(y|x)p(x)$ → X y

Parameters → Program → Observations

Frank Wood, Jan-Willem Van De Meent & Vikash Mansinghka. A New Approach to Probabilistic Programming Inference. AISTATS 2014, Reykjavik, JMLR, 1024-1032

---

Define model → Choose inference method → Derive algorithm by hand → Implement algorithm (e.g. Matlab) → Revise model/method → Re-implement algorithm (e.g. C++/C#)

Define model → Write model as *probabilistic program* → Apply inference engine → Revise model/engine settings

Image credit to John WINN (2010)

---

Image credit to Frank Wood (2016)

Graphical Models: BUGS, STAN
Factor Graphs: Factorie, Infer.NET

PL | AI | ML | STATS

2010 — Figaro, HANSAI, WebPPL, Venture, Anglican, Probabilistic-C, LibBi, STAN

Λₒ, Church, Infer.NET, Factorie, JAGS, ProbLog, Blog

2000 — IBAL, Prism, KMP, WinBUGS

1990 — BUGS, Simula, ALGOL 60, Prolog
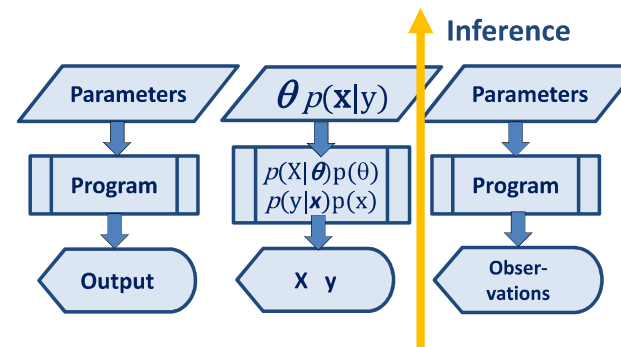
- 1940s: connecting wires to represent 0s and 1s
- 1950s: assemblers, FORTRAN, COBOL, LISP
- 1960s: ALGOL, BCPL(-> B -> C], SIMULA
- 1970s: Prolog, FP, ML, Miranda
- 1980s: Eiffel, C++
- 1990s: Haskell, Java, Python

---

- https://github.com/pymc-devs/pymc
- http://infernet.azurewebsites.net/
- http://mc-stan.org/
- https://github.com/p2t2/figaro
- https://sites.google.com/site/bloginference/
- http://projects.csail.mit.edu/church/wiki/Church
- http://factorie.cs.umass.edu/
- http://www.openbugs.net/w/FrontPage
- http://mcmc-jags.sourceforge.net/

---

Daniel Ritchie, Paul Horsfall & Noah D Goodman 2016. Deep Amortized Inference for Probabilistic Programs. arXiv:1610.05735.

Diederik P Kingma & Max Welling 2013. Auto-encoding variational Bayes. arXiv:1312.6114 (1983 citations as of 13.05.2018 07:00)
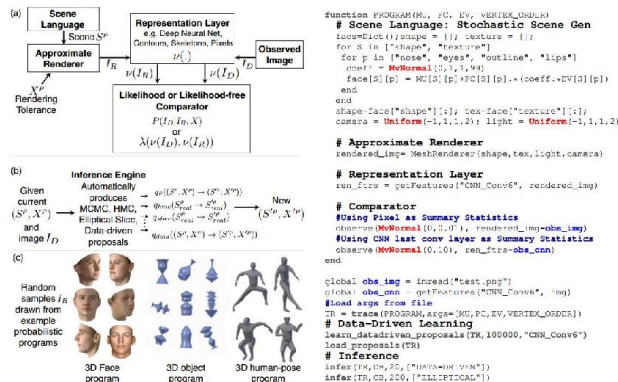
**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

$\theta, \phi \leftarrow$ Initialize parameters
**repeat**
 $\mathbf{X}^M \leftarrow$ Random minibatch of $M$ datapoints (drawn from full dataset)
 $\epsilon \leftarrow$ Random samples from noise distribution $p(\epsilon)$
 $\mathbf{g} \leftarrow \nabla_{\theta,\phi} \widetilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$ (Gradients of minibatch estimator (8))
 $\theta, \phi \leftarrow$ Update parameters using gradients g (e.g. SGD or Adagrad [DHS10])
**until** convergence of parameters $(\theta, \phi)$
**return** $\theta, \phi$

---

v1 [cs.AI] 17 Apr 2018

### Deep Probabilistic Programming Languages: A Qualitative Study

Guillaume Baudart
IBM Research
guillaume.baudart@ibm.com

Martin Hirzel
IBM Research
hirzel@us.ibm.com

Louis Mandel
IBM Research
lmandel@us.ibm.com

**ABSTRACT**

Deep probabilistic programming languages try to combine the advantages of deep learning with those of probabilistic programming languages. If successful, this would be a big step forward in machine learning and programming languages. Unfortunately, as of now, this new crop of languages is hard to use and understand. This paper addresses this problem directly by explaining deep probabilistic programming languages and indirectly by characterizing their current strengths and weaknesses.

**CCS CONCEPTS**
•Theory of computation → Probabilistic computation; •Computing methodologies → Neural networks; •Software and its engineering → Domain specific languages;

**KEYWORDS**
DL, PPL, DSL

**1 INTRODUCTION**

A deep probabilistic programming language (PPL) is a language for specifying both deep neural networks and probabilistic models. In other words, a deep PPL draws upon programming languages,

These frameworks provide automatic differentiation (users need not manually calculate gradients for gradient descent), GPU support (to efficiently execute vectorized computations), and Python-based embedded domain-specific languages [18].

Deep PPLs, which have emerged just recently [29–32], aim to combine the benefits of PPLs and DL. Ideally, programs in deep PPLs would overtly represent uncertainty, yield explainable models, and require only a small amount of training data; be easy to write in a well-designed programming language; and match the breakthrough accuracy and fast training times of DL. Realizing all of these promises would yield tremendous advantages. Unfortunately, this is hard to achieve. Some of the strengths of PPLs and DL are seemingly at odds, such as explainability vs. automated feature engineering, or learning from small data vs. optimizing for large data. Furthermore, the barrier to entry for work in deep PPLs is high, since it requires non-trivial background in fields as diverse as statistics, programming languages, and deep learning. To tackle this problem, this paper characterizes deep PPLs, thus lowering the barrier to entry, providing a programming-languages perspective early on when it can make a difference, and shining a light on gaps that the community should try to address.

This paper uses the Stan PPL as representative of the state of the art in regular (not deep) PPLs [9]. Stan is a main-stream, mature,

---

**Supervised learning**
x examples and labels
Θ parameters, boundaries

**Unsupervised learning**
x examples
Θ latent structures

**Reinforcement learning**
x actions, observations, rewards
Θ learned policy & world model

**Deep learning**
x training examples
Θ network weights    rg

Scott Cheng-Hsin Yang & Patrick Shafto 2017. Explainable Artificial Intelligence via Bayesian Teaching. NIPS 2017 Workshop Machine Teaching. Long Beach (CA).

---

| x | y |
|---|---|
| program source code | program output |
| scene description | image |
| policy and world | rewards |
| cognitive process | behavior |
| simulation | constraint |

Image credit to Frank Wood (2016)

---

Vikash K. Mansinghka, Tejas D. Kulkarni, Yura N. Perov & Josh Tenenbaum. Approximate Bayesian image interpretation using generative probabilistic graphics programs. In: Burges, Christopher J. C., Bottou, Leon, Welling, Max, Ghahramani, Zhoubin & Weinberger, Kilian Q., eds. Advances in Neural Information Processing Systems, 2013 Lake Tahoe. NIPS, 1520-1528.

### Approximate Bayesian Image Interpretation using Generative Probabilistic Graphics Programs

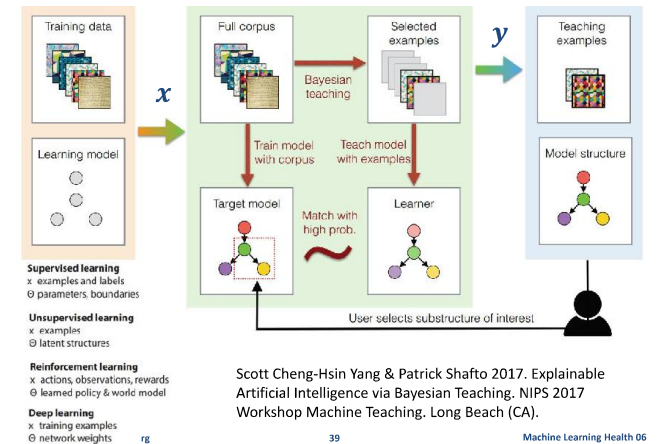Vikash K. Mansinghka[1,2], Tejas D. Kulkarni[1,2], Yura N. Perov[1,2,3], and Joshua B. Tenenbaum[1,2]

[1]Computer Science and Artificial Intelligence Laboratory, MIT
[2]Department of Brain and Cognitive Sciences, MIT
[3]Institute of Mathematics and Computer Science, Siberian Federal University
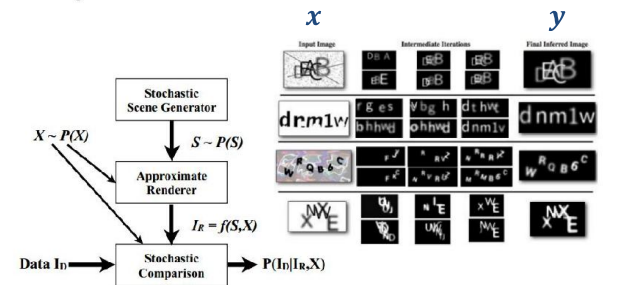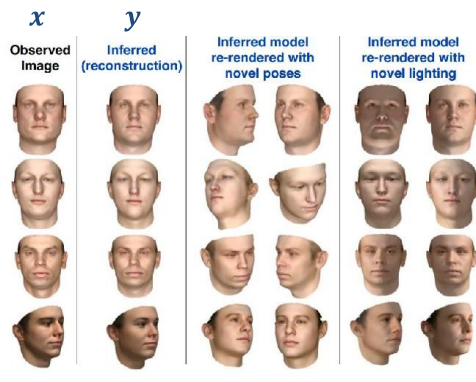
**Abstract**

The idea of computer vision as the Bayesian inverse problem to computer graphics has a long history and an appealing elegance, but it has proved difficult to directly implement. Instead, most vision tasks are approached via complex bottom-up processing pipelines. Here we show that it is possible to write short, simple probabilistic graphics programs that define flexible generative models and to automatically invert them to interpret real-world images. Generative probabilistic graphics programs (GPGP) consist of a stochastic scene generator, a renderer based on graphics software, a stochastic likelihood model linking the renderer's output and the data, and latent variables that adjust the fidelity of the renderer and the tolerance of the likelihood. Representations and algorithms from computer graphics are used as the deterministic backbone for highly approximate and stochastic generative models. This formulation combines probabilistic programming, computer graphics, and approximate Bayesian computation, and depends only on general-purpose, automatic inference techniques. We describe two applications: reading sequences of degraded and adversarially obscured characters, and inferring 3D road models from vehicle-mounted camera images. Each of the probabilistic graphics programs we present relies on under 20 lines of probabilistic code, and

---

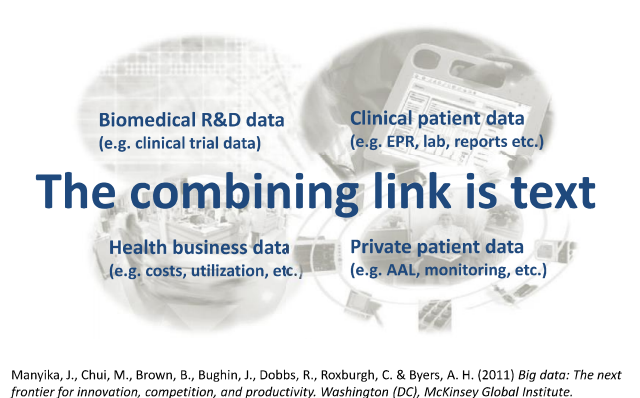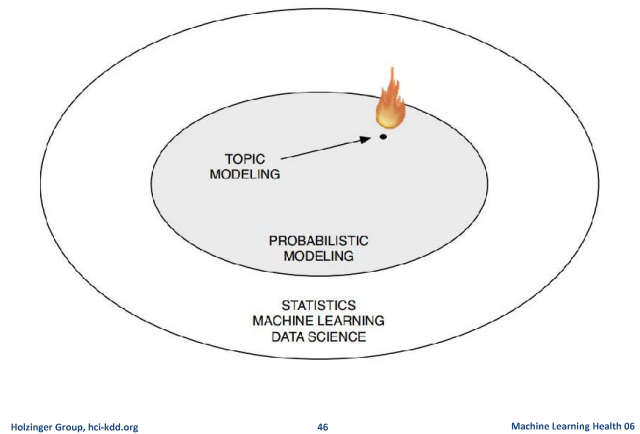$$P(S|I_D) \propto \int P(S)P(X)\delta_{f(S,X)}(I_R)P(I_D|I_R,X)dX$$



Vikash K. Mansinghka, Tejas D. Kulkarni, Yura N. Perov & Josh Tenenbaum. Approximate Bayesian image interpretation using generative probabilistic graphics programs. In: Burges, Christopher J. C., Bottou, Leon, Welling, Max, Ghahramani, Zhoubin & Weinberger, Kilian Q., eds. Advances in Neural Information Processing Systems, 2013 Lake Tahoe. NIPS, 1520-1528.

---

Kulkarni, Kohli, Tenenbaum & Mansinghka. Picture: A probabilistic programming language for scene perception. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 4390-4399.
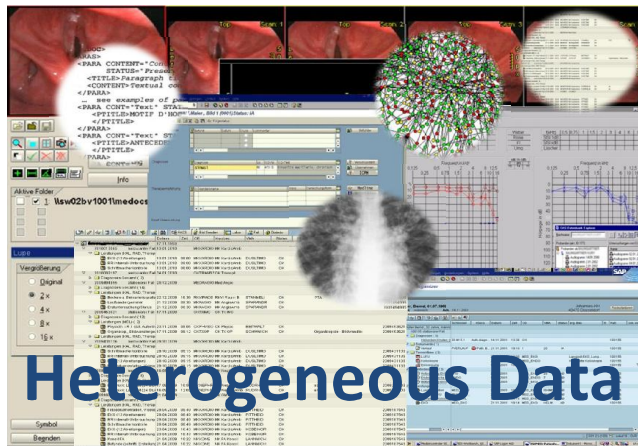
---

x    y



Observed Image | Inferred (reconstruction) | Inferred model re-rendered with novel poses | Inferred model re-rendered with novel lighting

Kulkarni, Kohli, Tenenbaum & Mansinghka. Picture: A probabilistic programming language for scene perception. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 4390-4399.

---



## 03 Probabilistic Topic Models

## Slide 46

TOPIC MODELING

PROBABILISTIC MODELING

STATISTICS
MACHINE LEARNING
DATA SCIENCE

## Slide 47



# Heterogeneous Data

## Slide 48



## Slide 49

**Biomedical R&D data** (e.g. clinical trial data)

**Clinical patient data** (e.g. EPR, lab, reports etc.)

# The combining link is text

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

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

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

## Slide 50

## Slide 51

Gerard M. Salton, Andrew Wong & Chungshu S. Yang 1975. Vector-Space Model for automatic indexing. Communications of the ACM, 18, (11), 613-620, doi:10.1145/361219.361220.

## Slide 52



Gerard M. Salton, Andrew Wong & Chungshu S. Yang 1975. Vector-Space Model for automatic indexing. Communications of the ACM, 18, (11), 613-620, doi:10.1145/361219.361220.

## Slide 53

Example (1)    ⬡HCI-KDD⬡

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

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

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

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

## Slide 54

Example (2)    ⬡HCI-KDD⬡

$$D_{m \times n} = \begin{cases} w_{1,1} & w_{1,2} & \cdots & w_{1,n-1} & w_{1,n} \\ w_{2,1} & w_{2,2} & & w_{2,n-1} & w_{2,n} \\ \vdots & & \ddots & & \vdots \\ w_{m-1,1} & w_{m-1,2} & & w_{m-1,n} & w_{m-1,n} \\ w_{m,1} & w_{m,2} & \cdots & w_{m,n-1} & w_{m,n} \end{cases}$$

Example (3)　　HCI-KDD



leber

$d_j$

verdacht

hepatitis

---

leber

$d_j$

φ

verdacht

$q$

hepatitis

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

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 $\alpha$ and $\beta$, the joint distribution of a topic mixture $\theta$, a set of $N$ topics $z$, and a set of $N$ words $w$ is given by:

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

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

---

$$P(w|d_i) = \sum_{i=1}^{k} P(w|t_i) \cdot P(t_i|d_i)$$

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 |

*Columns sorte[d] probability word given to[pic]*

*D. Blei*

---

Proportions parameter · Per-word topic assignment · Topic parameter · Per-document topic proportions · Observed word · Topics

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

---

$$p(\beta, \theta, \mathbf{z} \mid \mathbf{w}) = \frac{p(\beta, \theta, \mathbf{z}, \mathbf{w})}{\int_\beta \int_\theta \sum_\mathbf{z} p(\beta, \theta, \mathbf{z}, \mathbf{w})}$$

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

---

MASSIVE DATA · GLOBAL HIDDEN STRUCTURE · Subsample data · Infer local structure · Update global structure

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

---

1: Initialize $\lambda^{(0)}$ randomly.
2: Set the step-size schedule $\rho_t$ appropriately.
3: **repeat**
4:   Sample a document $w_d$ uniformly from the data set.
5:   Initialize $\gamma_{dk} = 1$, for $k \in \{1, \dots, K\}$.
6:   **repeat**
7:     For $n \in \{1, \dots, N\}$ set

$$\phi_{dn}^k \propto \exp\{\mathbb{E}[\log \theta_{dk}] + \mathbb{E}[\log \beta_{k,w_{dn}}]\}, \; k \in \{1, \dots, K\}.$$

8:     Set $\gamma_d = \alpha + \sum_n \phi_{dn}$.
9:   **until** local parameters $\phi_{dn}$ and $\gamma_d$ converge.
10:   For $k \in \{1, \dots, K\}$ set intermediate topics

$$\hat{\lambda}_k = \eta + D \sum_{n=1}^N \phi_{dn}^k w_{dn}.$$

11:   Set $\lambda^{(t)} = (1 - \rho_t)\lambda^{(t-1)} + \rho_t \hat{\lambda}$.
12: **until** forever

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

---

1. Sample a document
2. Estimate the local variational parameters using the current topics
3. Form intermediate topics from those local parameters
4. 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.

---

KNOWLEDGE · DATA · Make assumptions · Discover patterns · Predict & Explore

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

---

REUSABLE VARIATIONAL FAMILIES · MASSIVE DATA · ANY MODEL · BLACK BOX VARIATIONAL INFERENCE · $p(\beta, \mathbf{z} \mid \mathbf{x})$

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.

# 04 Knowledge Representation in Network Medicine

---

Image credit to Anna Goldenberg, Toronto

---

Pleotropic effects

Epistatic effects

Image credit to Eric Xing, Carnegie Mellon University, Pittsburgh

---

Nature Reviews | Molecular Cell Biology

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

---

Image credit to Anna Goldenberg, Toronto

---

# Nodes: 641
# Edges: 1250

Agent
Condition
Pharmacological Group
Other Documents

Average Degree: 3.888
Average Path Length: 4.683
Network Diameter: 9

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)

---

- Nodes
  - drugs
  - clinical guidelines
  - patient conditions (indication, contraindication)
  - pharmacological groups
  - tables and calculations of medical scores
  - algorithms and other medical documents
- Edges: 3 crucial types of relations inducing medical relevance between two active substances
  - pharmacological groups
  - indications
  - contra-indications

---

## Example for finding related structures

Relationship between
Adrenaline (center black node) and
Dobutamine (top left black node)
Blue: Pharmacological Group
Dark red: Contraindication;
Light red: Condition

Green nodes (from dark to light):
1. Application (one or more indications + corresponding dosages)
2. Single indication with additional details (e. g. "VF after 3$^{rd}$ Shock")
3. Condition (e.g. VF, Ventricular Fibrillation)

---

## Interactive Visual Data Mining

http://ophid.utoronto.ca/navigator



JURISICA LAB
IBM Life Sciences Discovery Center

HCI-KDD

Otasek, D., Pastrello, C., Holzinger, A. & Jurisica, I. 2014. Visual Data Mining: Effective Exploration of the Biological Universe. In: Holzinger, A. & Jurisica, I. (eds.) Interactive Knowledge Discovery and Data Mining in Biomedical Informatics: State-of-the-Art and Future Challenges. Lecture Notes in Computer Science LNCS 8401. Heidelberg, Berlin: Springer, pp. 19–34, doi:10.1007/978-3-662-43968-5_2.

---



---

## Example: Graph Entropy Measures



▪ Engineering
▪ Computer Science
▨ Physics
▥ Humanities
□ unkown

Holzinger et. al. 2013. On Graph Entropy Measures for Knowledge Discovery from Publication Network Data. In: LNCS 8127, 354-362.

---

## Some selected open problems

- **Problem:** What is the max. number of edges of an Relative Neighborhood Graph in R3 ? No supra-linear lower bound is known.
- **Problem:** What is the structural interpretation of graph measures ? They are mappings which maps graphs to the reals. Thus, they can be understood as graph complexity measures and investigating their structural interpretation relates to understand what kind of structural complexity they detect.
- **Problem:** It is important to visualize large networks meaningfully. So far, there has been a lack of interest to develop efficient software beyond the available commercial software.
- **Problem:** Are multi-touch interaction graphs structurally similar to other graphs (from known graph classes)? This calls for a comparison of graph classes and their structural characteristics.
- **Problem:** Which graph measures are suitable to determine the complexity of multi-touch interaction graphs? Does this lead to any meaningful classification based on their topology?
- **Problem:** What is interesting? Where to start the interaction?

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

---

## Example: The brain is a complex network



Van Den Heuvel, M. P. & Hulshoff Pol, H. E. (2010) Exploring the brain network: a review on resting-state fMRI functional connectivity. *European Neuropsycho-pharmacology, 20, 8, 519-534.*

---

## Representative Examples of disease complexes

Examples of 4 functional networks driving the development of different anatomical structures in the human heart of a 37-day old human embryo



Lage, K. et. al (2010) Dissecting spatio-temporal protein networks driving human heart development and related disorders. *Molecular systems biology, 6, 1, 1-9.*

---

## Example: Cell-based therapy



Lage et. al (2010)

---

## Identifying Networks in Disease Research



Schadt, E. E. & Lum, P. Y. (2006) Reverse engineering gene networks to identify key drivers of complex disease phenotypes. *Journal of lipid research, 47, 12, 2601-2613.*

## Three main types of biomedical networks

Transcriptional regulatory network with two components:
TF = transcription factor
TG = target genes
(TF regulates the transcription of TG)

Protein-Protein interaction network

Metabolic network (constructed considering the reactants, chemical reactions and enzymes)

Costa, L. F., Rodrigues, F. A. & Cristino, A. S. (2008) Complex networks: the key to systems biology. *Genetics and Molecular Biology, 31, 3, 591–601.*

---

## Example Transcriptional Regulatory Network



Salgado, H., Santos-Zavaleta, A., Gama-Castro, S., Peralta-Gil, M., Peñaloza-Spínola, M. I., Martínez-Antonio, A., Karp, P. D. & Collado-Vides, J. 2006. The comprehensive updated regulatory network of Escherichia coli K-12. *BMC bioinformatics, 7, (1), 5.*

---

## Network Representations of Protein Complexes



A — Protein complex
B — True PPI topology
C — Matrix-Model
D — Spoke-Model

Wang, Z. & Zhang, J. Z. (2007) In search of the biological significance of modular structures in protein networks. PLoS Computational Biology, 3, 6, 1011-1021.

---

## Correlated Motif Mining (CMM)



$$V_X \cap V_Y$$

$V_X$    $V_Y$

Boyen, P., Van Dyck, D., Neven, F., van Ham, R. C. H. J. & van Dijk, A. (2011) SLIDER: A Generic Metaheuristic for the Discovery of Correlated Motifs in Protein-Protein Interaction Networks. *Computational Biology and Bioinformatics, IEEE/ACM Transactions on, 8, 5, 1344-1357.*

---

## Steepest Ascent Algorithm applied to CMM

**Input:** PPI-network $G = (V, E, \lambda)$, $\ell, d \in \mathbb{N}$, $d < \ell$
**Output:** $\{X^*, Y^*\}$ best correlated motif pair found in $G$
1: $\{X^*, Y^*\} \leftarrow \text{randomMotifPair}()$
2: $maxsup \leftarrow f(\{X^*, Y^*\}, G)$
3: $sup \leftarrow -\infty$
4: **while** $maxsup > sup$ **do**
5: $\quad \{X, Y\} \leftarrow \{X^*, Y^*\}$
6: $\quad sup \leftarrow maxsup$
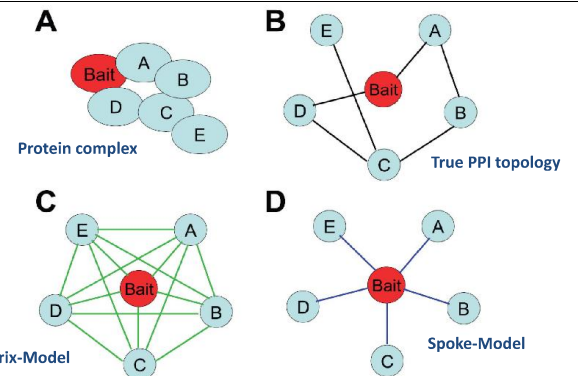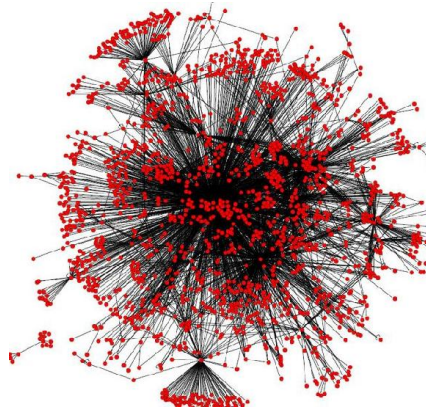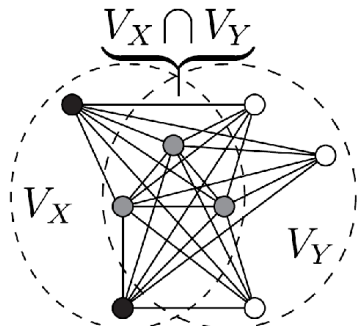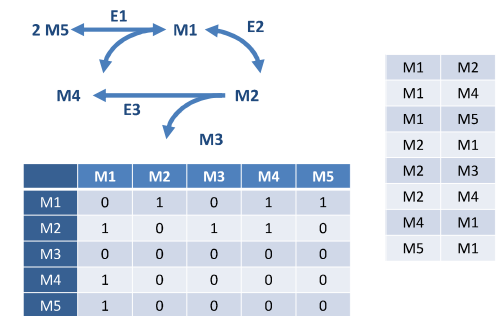7: $\quad$ **for all** $\{X', Y'\} \in N(\{X, Y\})$ **do**
8: $\quad\quad$ **if** $f(\{X', Y'\}, G) > maxsup$ **then**
9: $\quad\quad\quad \{X^*, Y^*\} \leftarrow \{X', Y'\}$
10: $\quad\quad\quad maxsup \leftarrow f(\{X', Y'\}, G)$

Boyen et al. (2011)

---

## Metabolic Network



|     | M1 | M2 | M3 | M4 | M5 |
|-----|----|----|----|----|----|
| M1  | 0  | 1  | 0  | 1  | 1  |
| M2  | 1  | 0  | 1  | 1  | 0  |
| M3  | 0  | 0  | 0  | 0  | 0  |
| M4  | 1  | 0  | 0  | 0  | 0  |
| M5  | 1  | 0  | 0  | 0  | 0  |

| M1 | M2 |
|----|----|
| M1 | M4 |
| M1 | M5 |
| M2 | M1 |
| M2 | M3 |
| M2 | M4 |
| M4 | M1 |
| M5 | M1 |

**Matrix contains many sparse elements - In this case it is computationally more efficient to represent the graph as an adjacency list**

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

---

## Metabolic networks are usually big … big data ☺



Schmid, A. K., Reiss, D. J., Pan, M., Koide, T. & Baliga, N. S. (2009) A single transcription factor regulates evolutionarily diverse but functionally linked metabolic pathways in response to nutrient availability. *Molecular Systems Biology, 5, 1-9.*

http://www.nature.com/msb/journal/v5/n1/fig_tab/msb200940_F6.html

---

## Using EPRs to Discover Disease Correlations



*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., Søeby, 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.*

---

## Heatmap of disease-disease correlations (ICD)



Roque, F. S. et al (2011) Using Electronic Patient Records to Discover Disease Correlations and Stratify Patient Cohorts. *PLoS Comput Biol, 7, 8, e1002141.*

**Example: ὁμολογέω (homologeo)**

T0499    T0498

He, Y., Chen, Y., Alexander, P., Bryan, P. N. & Orban, J. (2008) NMR structures of two designed proteins with high sequence identity but different fold and function. Proceedings of the National Academy of Sciences, 105, 38, 14412.

Tyr 45
Ala 20
Ile 30
Phe 30
Leu 45
Leu 20

T0499  TTYKLILNLKQAKEEAIKELVDAGTAEKYFKLIANAKTVEGVWTYKDEIKTFTVTE

T0498  TTYKLILNLKQAKEEAIKELVDAGTAEKYIKLIANAKTVEGVWTLKDEIKTFTVTE

---

**Conclusion**

- 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** underline{personalized and molecular medicine} **in the future.**

---

# 05 Machine Learning on Graphs Relevant for Health Informatics

---

**Example: Lymphoma is the most common blood cancer**

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 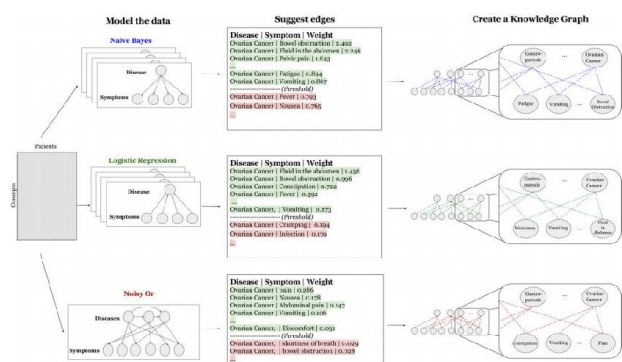of white blood cell, grow and multiply uncontrollably. Cancerous lymphocytes can travel to many parts of the body, including the lymph nodes, spleen, bone marrow, blood, or other organs, and form a mass called a tumor. The body has two main types of lymphocytes that can develop into lymphomas: B-lymphocytes (B-cells) and T-lymphocytes (T-cells).
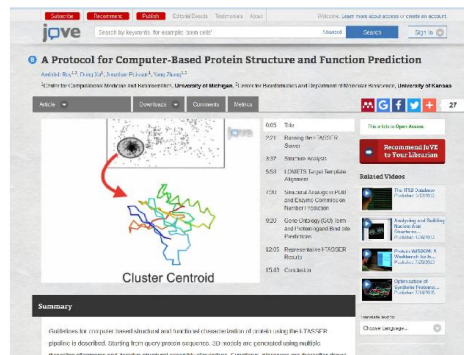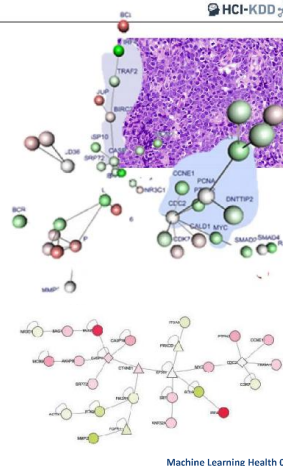
www.lymphoma.org

http://imagebank.hematology.org/

---

**ML tasks on graphs**

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

Dittrich, M. T., Klau, G. W., Rosenwald, A., Dandekar, T. & Müller, T. 2008. Identifying functional modules in protein–protein interaction networks: an integrated exact approach. Bioinformatics, 24, (13), i223-i231.

---

**SCIENTIFIC REPORTS**

OPEN **Learning a Health Knowledge Graph from Electronic Medical Records**

Received: 3 March 2017
Accepted: 1 June 2017
Published online: 20 July 2017

Maya Rotmensch[1], Yoni Halpern[2], Abdulhakim Tlimat[3], Steven Horng[3,4] & David Sontag[5,6]

Demand for clinical decision support systems in medicine and self-diagnostic symptom checkers has substantially increased in recent years. Existing platforms rely on knowledge bases manually compiled through a labor-intensive process or automatically derived using simple pairwise statistics. This study explored an automated process to learn high quality knowledge bases linking diseases and symptoms directly from electronic medical records. Medical concepts were extracted from 273,174 de-identified patient records and maximum likelihood estimation of three probabilistic models was used to automatically construct knowledge graphs: logistic regression, naive Bayes classifier and a Bayesian network using noisy OR gates. A graph of disease-symptom relationships was elicited from the learned parameters and the constructed knowledge graphs were evaluated and validated, with permission, against Google's manually-constructed knowledge graph and against expert physician opinions. Our study shows that direct and automated construction of high quality health knowledge graphs from medical records using rudimentary concept extraction is feasible. The noisy OR model produces a high quality knowledge graph reaching precision of 0.85 for a recall of 0.6 in the clinical evaluation. Noisy OR significantly outperforms all tested models across evaluation frameworks ($p < 0.01$).

---

**Workflow for modeling relationship disease-symptom**

Model the data | Suggest edges | Create a Knowledge Graph

Maya Rotmensch, Yoni Halpern, Abdulhakim Tlimat, Steven Horng & David Sontag 2017. Learning a Health Knowledge Graph from Electronic Medical Records. Scientific Reports, 7, 5994, doi:10.1038/s41598-017-05778-z.

---

**From structure to function**

A Protocol for Computer-Based Protein Structure and Function Prediction

Cluster Centroid

http://www.jove.com/video/3259/a-protocol-for-computer-based-protein-structure-function

---

**Interesting: Hubs tend to link to small degree nodes**

Nodes: proteins
Links:   physical interactions (binding)

Puzzling pattern:
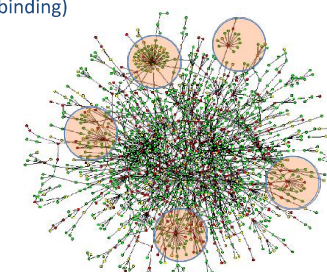Hubs tend to link to small degree nodes.
Why is this puzzling?
In a random network, the probability that a node with degree $k$ links to a node with degree $k'$ is:
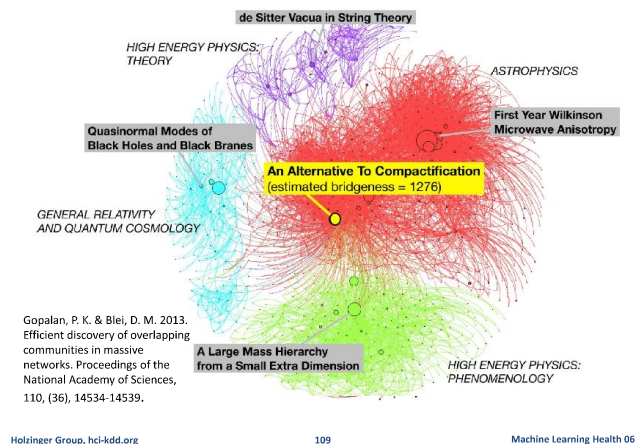
$$p_{kk'} = \frac{kk'}{2L}$$

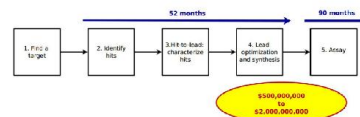$k$=50, $k'$=13, N=1,458, L=1746

$p_{50,13} = 0.15$      $p_{2,1} = 0.0004$

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

de Sitter Vacua in String Theory

HIGH ENERGY PHYSICS: THEORY

ASTROPHYSICS

First Year Wilkinson Microwave Anisotropy

Quasinormal Modes of Black Holes and Black Branes

An Alternative To Compactification (estimated bridgeness = 1276)

GENERAL RELATIVITY AND QUANTUM COSMOLOGY

A Large Mass Hierarchy from a Small Extra Dimension

HIGH ENERGY PHYSICS: PHENOMENOLOGY

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?
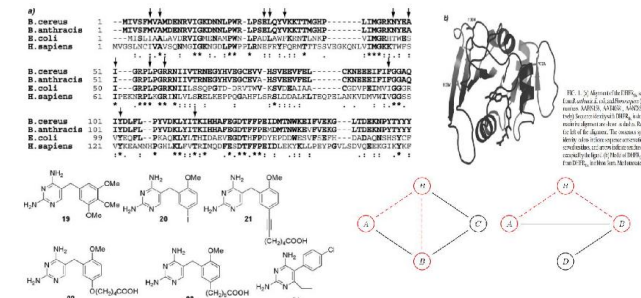- How can we learn the parameters and structure of a graphical model?



The chemical space

- $10^{60}$ possible small or-ganic molecules
- $10^{22}$ stars in the observ-able 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

---

# Thank you!

---

# Questions

---

- 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?

---

- 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?
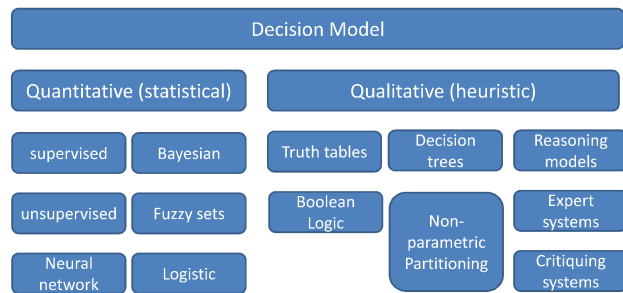- Describe the LDA-model and its application for topic modelling!

---

- 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 trail – 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?

- 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 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.

# Appendix

# 1) Reasoning under Uncertainty

## Remember: Taxonomy of Decision Support Models

Decision Model

- Quantitative (statistical)
  - supervised
  - Bayesian
  - unsupervised
  - Fuzzy sets
  - Neural network
  - Logistic
- Qualitative (heuristic)
  - Truth tables
  - Decision trees
  - Reasoning models
  - Boolean Logic
  - Non-parametric Partitioning
  - Expert systems
  - Critiquing systems

Bemmel, J. H. v. & Musen, M. A. (1997) *Handbook of Medical Informatics. Heidelberg, Springer.*
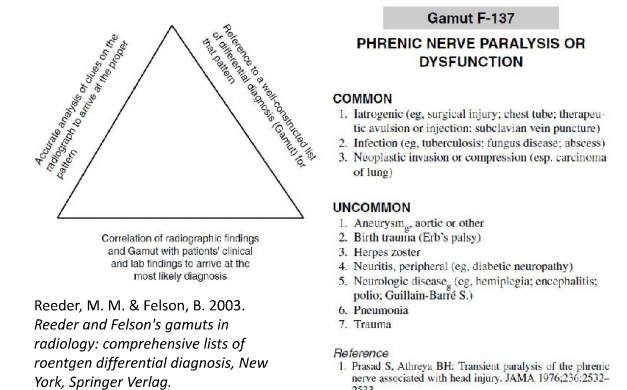
## Dealing with uncertainty in the real world

- 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 – 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

## 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.*

## MYCIN was *no* success in the clinical practice

https://www.youtube.com/watch?v=IVGWM0CKNWA ("real nurse triage")

## Gamuts: Triangulation to find diagnoses



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. Aneurysm, 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, Alireza BH: Transient paralysis of the phrenic nerve associated with head injury. JAMA 1976;236:2532–2533

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

**REEDER AND FELSON'S GAMUTS IN RADIOLOGY**

GAMUT G-25
EROSIVE GASTRITIS*

COMMON
1. Acute gastritis (eg, alcohol abuse)
2. Crohn's disease
3. Drugs (eg, aspirin, NSAID, steroids)
4. Helicobacter pylori infection
5. Idiopathic
6. [Normal areae gastricae]
7. Peptic ulcer; hyperacidity

UNCOMMON
1. Corrosive gastritis
2. Cryptosporidium antritis
3. [Lymphoma]
4. Opportunistic infection (eg, candidiasis (moniliasis); herpes simplex; cytomegalovirus)
5. Postoperative gastritis
6. Radiation therapy
7. Zollinger-Ellison S.; multiple endocrine neoplasia (MEN) S.

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

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

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

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

---

- 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

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

---

- **Type 1 Decisions:** related to the **diagnosis,** i.e. computers are used to assist in diagnosing a disease on the basis of the individual patient data. Questions include:
  - What is the probability that this patient has a myocardial infarction on the basis of given data (patient history, ECG, ...)?
  - What is the probability that this patient has acute appendices, given the signs and symptoms concerning abdominal pain?

- **Type 2 Decisions:** related to **therapy,** i.e. computers are used to select the best therapy on the basis of clinical evidence, e.g.:
  - What is the best therapy for patients of age $x$ and risks y, if an obstruction of more than $z$ % is seen in the left coronary artery?
  - What amount of insulin should be prescribed for a patient during the next 5 days, given the blood sugar levels and the amount of insulin taken during the recent weeks?

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

---

The future is in integrative ML, i.e. combining relational databases, ontologies and logic with probabilistic reasoning models and statistical learning – and algorithms that have good **scalability**

w Smokes(x) ∧ Friends(x,y) ⇒ Smokes(y)

Big data
Big models

Learns a model over 900,030,000 random variables

Van Den Broeck, G., Taghipour, N., Meert, W., Davis, J. & De Raedt, L. Lifted probabilistic inference by first-order knowledge compilation. Proceedings of the Twenty-Second international joint conference on Artificial Intelligence-Volume Volume Three, 2011. AAAI Press, 2178-2185.

---

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

h₁ = The identity of ORGANISM-1 is streptococcus
h₂ = PATIENT-1 is febrile
h₃ = The name of PATIENT-1 is John Jones

CF[h₁,E] = .8    : There is strongly suggestive evidence (.8) that the identity of ORGANISM-1 is streptococcus
CF[h₂,E] = –.3    : There is weakly suggestive evidence (.3) that PATIENT-1 is not febrile
CF[h₃,E] = +1    : It is definite (1) that the name of PATIENT-1 is John Jones

---

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

**PyMC3**

PyMC Pythonic Markov chain Monte Carlo

---

| 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 |

---

* Simple example: Nucleotide "A" may follow nucleotide "T" in the sequences more frequently for outcome X than for outcome Y,

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

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

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

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

$$P(\theta \mid D) = \frac{P(D \mid \theta) \cdot P(\theta)}{P(D)}$$

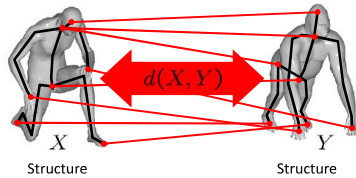Image Source: Dan Williams, Life Technologies, Austin TX

---

# 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.*

---



Rock    Hands    Scissors    Paper

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

---

## Similarity and Correspondence

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

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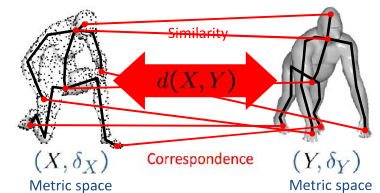


$d(X,Y)$

$X$     $Y$

Structure     Structure

Correspondence quality = structure similarity (distortion)

Minimum possible correspondence distortion

---

## Invariant Similarity



Similarity

$\tau$   $d(X,Y)$ transformation   $\sigma$

Invariant similarity

$d(\tau X, \sigma Y)$

Minimum possible correspondence $= d(X,Y)$

$\tau X$     $\sigma Y$

---

## Gromov-Hausdorff dist: finding the opt. correspondence

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



Michail Gromov (1943- )     Felix Hausdorff (1868-1942)

Similarity

$d(X,Y)$

$(X, \delta_X)$   Correspondence   $(Y, \delta_Y)$

Metric space     Metric space

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

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

**Discrete optimization over correspondences is NP hard !**

---

## Distinguish topological spaces

Counts the number of "i-dimensional holes"

$b_i$ is the "i-th Betti number"

**Enrico Betti (1823-1892)**     **Emmy Noether (1882-1935)**



$b_1=1$   $b_2=0$     $b_1=0$   $b_2=1$     $b_1=2$   $b_2=1$

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.*

---

## Structural Patterns are often hidden in weakly str. data

- 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 $f: M \longrightarrow \mathbb{R}^d$ 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 $(\Omega, \mathcal{F}, P)$, 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 $M$?

---

## Topological Data Mining



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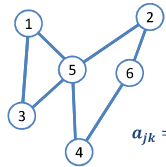
Holzinger, A. 2014. On Topological Data Mining. In: Lecture Notes in Computer Science, LNCS 8401. Berlin Heidelberg: Springer, pp. 331-356, doi:10.1007/978-3-662-43968-5_19.

# 06 Review of basic concepts, metrics and measures

## Complex Biological Systems key concepts

- 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

## Network Basics on the Example of Bioinformatics

$G(V, E)$ **Graph**
$V$ ...vertex
$E$ ...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.*

## Baby Stuff: Computational Graph Representation

Adjacency (ə-ˈjā-sᵊn(t)-sē) Matrix $A = (a_{jk})$

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

LEONHARD EULER 1707-1783

130

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

**Simple graph, symmetric, binary**

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

**Directed and weighted**

For more information: Diestel, R. (2010) *Graph Theory, 4th Edition. Berlin, Heidelberg, Springer.*
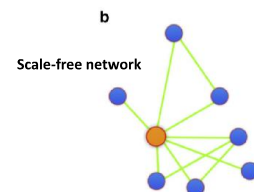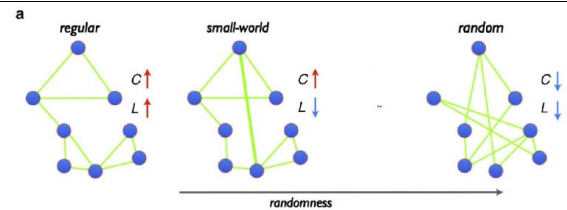
## Example: Tool for Node-Link Visualization

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

Fekete, J.-D. The infovis toolkit.  Information Visualization, INFOVIS 2004, 2004. IEEE, 167-174.
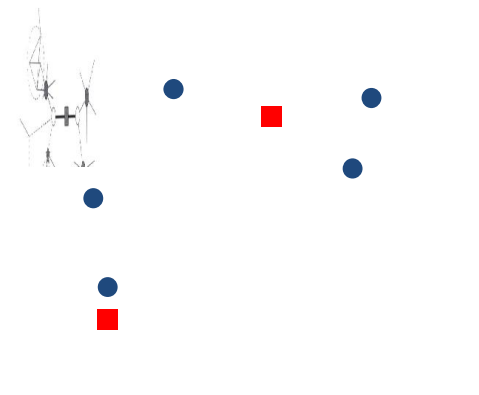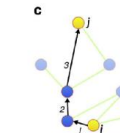
## Some Network Metrics (1/2)

**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:

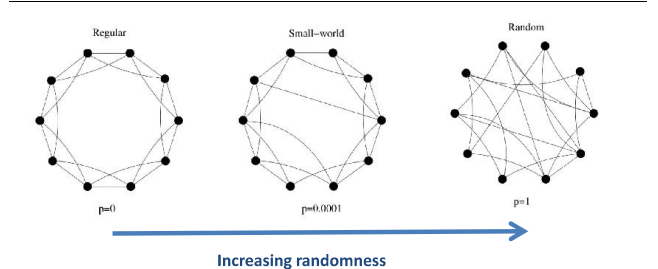$$C_i = \frac{2t_i}{k_i(k_i - 1)} \qquad C = \frac{1}{n} \sum_i C_i$$

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

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

Costa, L. F., Rodrigues, F. A., Travieso, G. & Boas, P. R. V. (2007) Characterization of complex networks: A survey of measurements. *Advances in Physics, 56, 1, 167-242.*
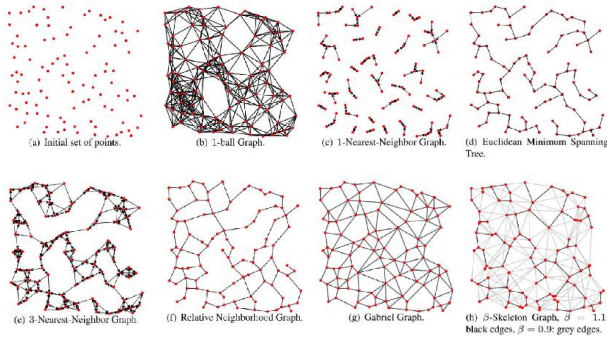
## Some Network Metrics (2/2)

- **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: $k_i = \sum_i a_{ij}$

**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).

## Network Topologies

a

regular          small-world          random

$C \uparrow$          $C \uparrow$          $C \downarrow$
$L \uparrow$          $L \downarrow$          $L \downarrow$

randomness

b

**Scale-free network**

Van Heuvel & Hulshoff (2010)

## Small-World Networks

Regular          Small-world          Random

p=0          p=0.0001          p=1

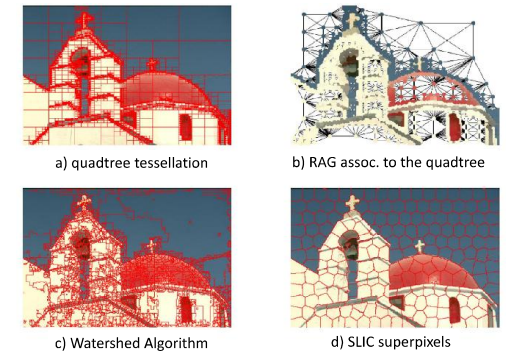**Increasing randomness**

29.000 citations ...

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) Initial set of points. (b) 1-ball Graph. (c) 1-Nearest-Neighbor Graph. (d) Euclidean Minimum Spanning Tree.

(e) 3-Nearest-Neighbor Graph. (f) Relative Neighborhood Graph. (g) Gabriel Graph. (h) $\beta$-Skeleton Graph, $\beta = 1.1$: black edges, $\beta = 0.9$: grey edges.
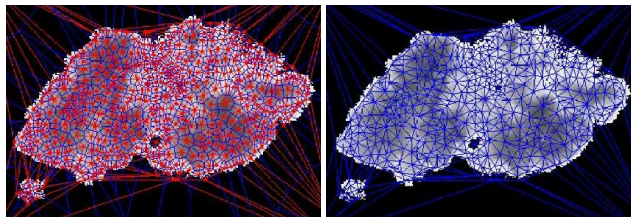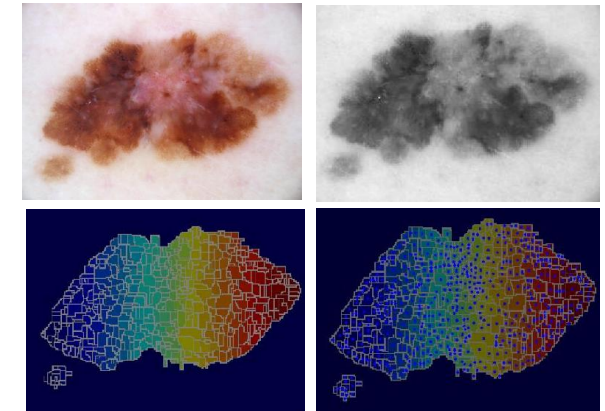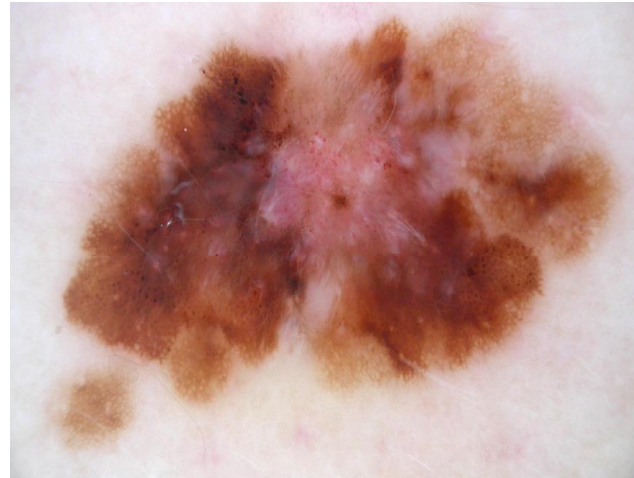
Lézoray, O. & Grady, L. 2012. Graph theory concepts and definitions used in image processing and analysis. In: Lézoray, O. & Grady, L. (eds.) Image Processing and Analysing With Graphs: Theory and Practice. Boca Raton (FL): CRC Press, pp. 1-24.

---

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

---

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.

---

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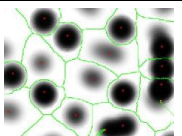
For Delaunay please refer to: Lee, D.-T. & Schachter, B. J. 1980. Two algorithms for constructing a Delaunay triangulation. Intl. Journal of Computer & Information Sciences, 9, (3), 219-242.

---

- 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
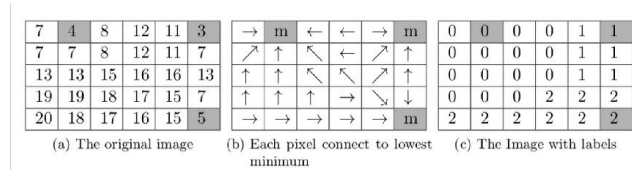
---

- 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.

- 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

---

(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

---

- 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**

---

- Study of complex networks started in the 1990s with the insight that real networks contain properties not present in random (Erdös-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).

---

EB    PB    TB

Personalized Medicine

Proteomics

Genomics

2003    2013    2023