

Andreas Holzinger

185.A83 Machine Learning for Health Informatics

2018S, VU, 2.0 h, 3.0 ECTS

Lecture 09 - Module 07 – Week 23 – 05.06.2018



## From Clinical Decision Support to explainable AI (ex-AI)

a.holzinger@hci-kdd.org

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



Holzinger Group hci-kdd.org

1

Machine Learning Health 08

- **DXplain** = a DSS from the Harvard Medical School, to assist making a diagnosis (clinical consultation), and also as an instructional instrument (education); provides a description of diseases, etiology, pathology, prognosis and up to 10 references for each disease;
- Etiology = in medicine (many) factors coming together to cause an illness (see causality)
- Explainable AI = Explainability = upcoming fundamental topic within recent AI; answering e.g. **why** a decision has been made
- **Expert-System** = emulates the decision making processes of a human expert to solve complex problems;
- **GAMUTS** in Radiology = Computer-Supported list of common/uncommon differential diagnoses;
- **ILIAD** = medical expert system, developed by the University of Utah, used as a teaching and testing tool for medical students in problem solving. Fields include Pediatrics, Internal Medicine, Oncology, Infectious Diseases, Gynecology, Pulmonology etc.
- Interpretability = there is no formal technical definition yet, but it is considered as a prerequisite for trust
- **MYCIN** = one of the early medical expert systems (Shortliffe (1970), Stanford) to identify bacteria causing severe infections, such as bacteremia and meningitis, and to recommend antibiotics, with the dosage adjusted for patient's body weight;
- **Reasoning** = cognitive (thought) processes involved in making medical decisions (clinical reasoning, medical problem solving, diagnostic reasoning);
- **Transparency** = opposite of opacity of black-box approaches, and connotes the ability to understand how a model works (that does not mean that it should always be understood, but that – in the case of necessity – it can be re-enacted)

Holzinger Group hci-kdd.org

4

Machine Learning Health 08



Holzinger Group hci-kdd.org

7

Machine Learning Health 08

- Decision support system (DSS)
- MYCIN – Rule Based Expert System
- GAMUTS in Radiology
- Reasoning under uncertainty
- Example: Radiotherapy planning
- Example: Case-Based Reasoning
- Explainable Artificial intelligence
- Re-trace > Understand > Explain
- Transparency > Trust > Acceptance
- Fairness > Transparency > Accountability
- Methods of Explainable AI

Holzinger Group hci-kdd.org

2

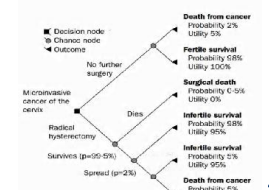
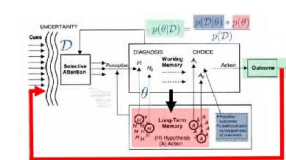
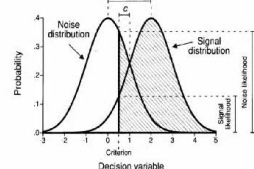
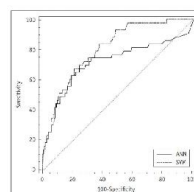
Machine Learning Health 08

- ... can apply your knowledge gained in the previous lectures to example systems of decision support;
- ... have an overview about the core principles and architecture of decision support systems;
- ... are familiar with the certainty factors as e.g. used in MYCIN;
- ... are aware of some design principles of DSS;
- ... have seen similarities between DSS and KDD on the example of computational methods in cancer detection;
- ... have seen basics of CBR systems;

Holzinger Group hci-kdd.org

5

Machine Learning Health 08



Holzinger Group hci-kdd.org

8

Machine Learning Health 08

- **Case-based reasoning (CBR)** = process of solving new problems based on the solutions of similar past problems;
- **Certainty factor model (CF)** = a method for managing uncertainty in rule-based systems;
- **CLARION** = Connectionist Learning with Adaptive Rule Induction ON-line (CLARION) is a cognitive architecture that incorporates the distinction between implicit and explicit processes and focuses on capturing the interaction between these two types of processes. By focusing on this distinction, CLARION has been used to simulate several tasks in cognitive psychology and social psychology. CLARION has also been used to implement intelligent systems in artificial intelligence applications.
- **Clinical decision support (CDS)** = process for enhancing health-related decisions and actions with pertinent, organized clinical knowledge and patient information to improve health delivery;
- **Clinical Decision Support System (CDSS)** = expert system that provides support to certain reasoning tasks, in the context of a clinical decision;
- **Collective Intelligence** = shared group (symbolic) intelligence, emerging from cooperation/competition of many individuals, e.g. for consensus decision making;
- **Crowdsourcing** = a combination of "crowd" and "outsourcing" coined by Jeff Howe (2006), and describes a distributed problem-solving model; example for crowdsourcing is a public software beta-test;
- **Decision Making** = central cognitive process in every medical activity, resulting in the selection of a final choice of action out of several alternatives;
- **Decision Support System (DSS)** = is an IS including knowledge based systems to interactively support decision-making activities, i.e. making data useful;

Holzinger Group hci-kdd.org

3

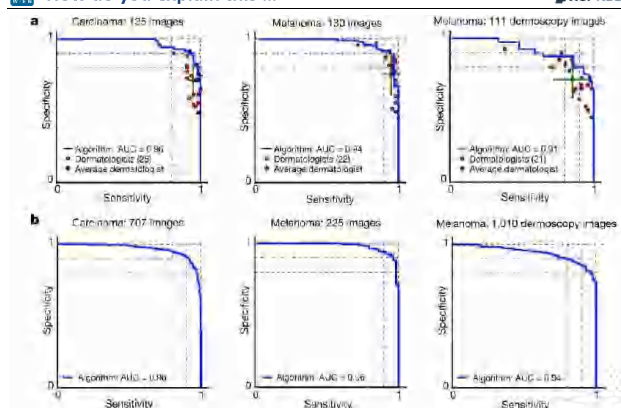
Machine Learning Health 08

- **00 Reflection – follow-up from last lecture**
- **01 Decision Support Systems (DSS)**
- **02 Computers help making better decisions?**
- **03 History of DSS = History of AI**
- **04 Example: Towards Personalized Medicine**
- **05 Example: Case Based Reasoning (CBR)**
- **06 Towards Explainable AI**
- **07 Some methods of explainable AI**

Holzinger Group hci-kdd.org

6

Machine Learning Health 08



Andre Esteve, Brett Kuperl, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau & Sebastian Thrun 2017. Dermatologist-level classification of skin cancer with deep neural networks. Nature, 542, (7639), 115-118

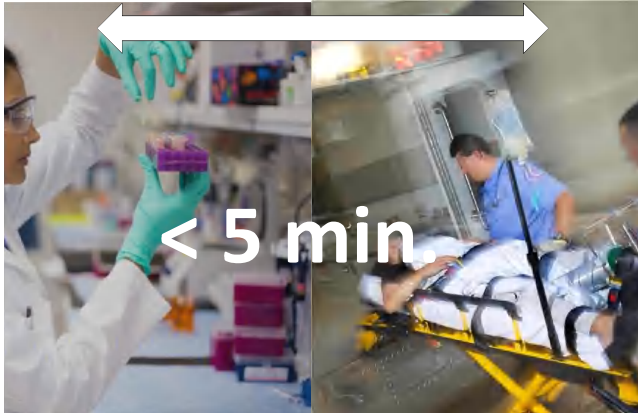
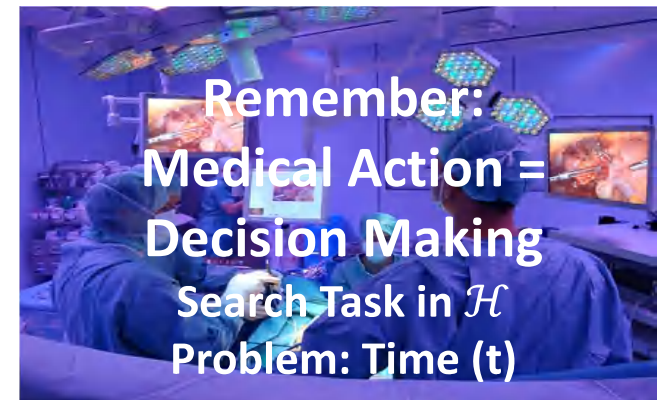
Holzinger Group hci-kdd.org

9

Machine Learning Health 08

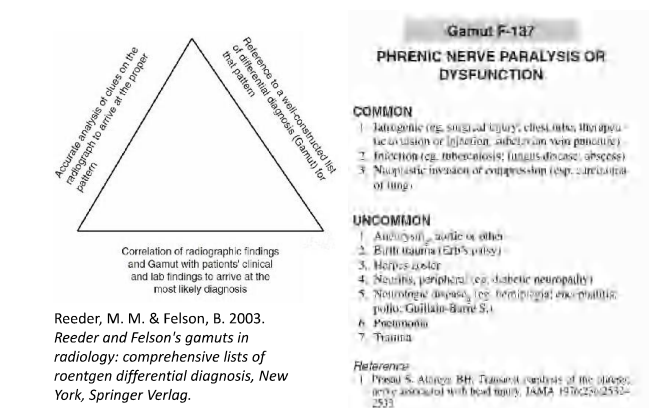
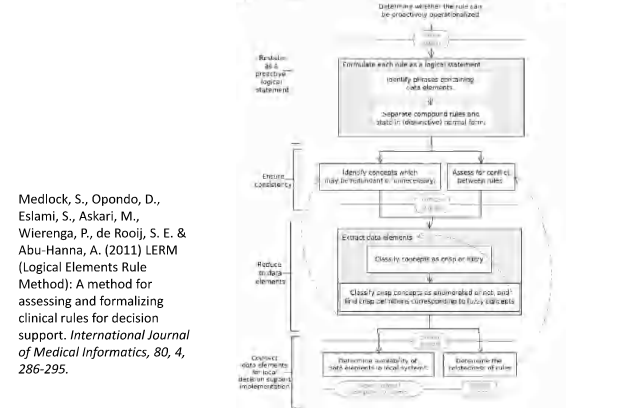
- Remember: Medicine is an complex application domain – dealing most of the time with **probable information!**
- Some challenges include:
  - (a) defining hospital system architectures in terms of generic tasks such as diagnosis, therapy planning and monitoring to be executed for (b) medical reasoning in (a);
  - (c) patient information management with (d) minimum uncertainty.
- Other challenges include: (e) knowledge acquisition and encoding, (f) human-computer interface and interaction; and (g) system integration into existing clinical legacy and proprietary environments, e.g. the enterprise hospital information system; to mention only a few.

## 01 Decision Support Systems

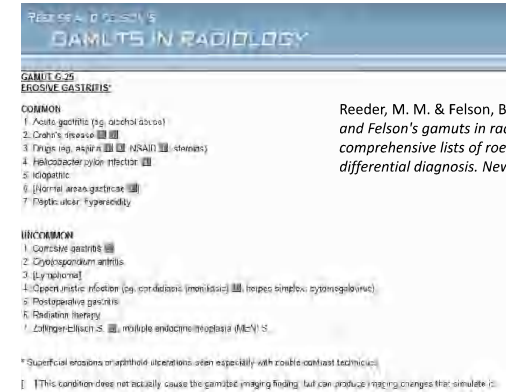


- 400 BC Hippocrates (460-370 BC), father of western medicine:
  - A medical record should accurately reflect the course of a disease
  - A medical record should indicate the probable cause of a disease
- 1890 William Osler (1849-1919), father of modern western medicine
  - Medicine is a science of uncertainty and an art of probabilistic decision making
- Today
  - Prediction models are based on data features, patient health status is modelled as high-dimensional feature vectors ...

- Clinical guidelines are **systematically** developed documents to assist doctors and patient decisions about appropriate care;
- In order to build DS, based on a guideline, it is **formalized** (transformed from natural language to a logical algorithm), and
- implemented** (using the algorithm to program a DSS);
- To increase the quality of care, they must be linked to a process of care, for example:
  - "80% of diabetic patients should have an HbA1c below 7.0" could be linked to processes such as:
    - "All diabetic patients should have an annual HbA1c test" and
    - "Patients with values over 7.0 should be rechecked within 2 months."
- Condition-action rules** specify one or a few conditions which are linked to a specific action, in contrast to narrative guidelines which describe a series of branching or iterative decisions unfolding over time.
- Narrative guidelines and clinical rules are two ends of a continuum of clinical care standards.

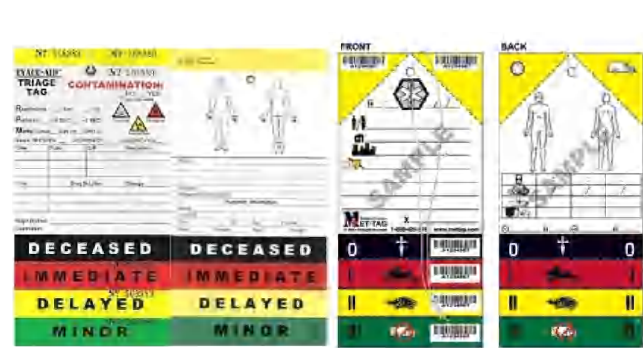






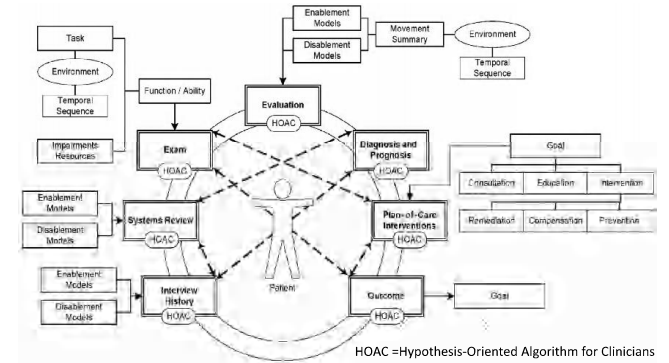
Reeder, M. M. & Felson, B. (2003) *Reeder and Felson's gamuts in radiology: comprehensive lists of roentgen differential diagnosis*. New York, Springer

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



Iseron, K. V. & Moskop, J. C. 2007. Triage in Medicine, Part I: Concept, History, and Types. *Annals of Emergency Medicine*, 49, (3), 275-281.

Image Source: <http://store.zomedtech.com>



Schenkman, M., Deutsch, J. E. & Gill-Body, K. M. (2006) An Integrated Framework for Decision Making in Neurologic Physical Therapist Practice. *Physical Therapy*, 86, 12, 1681-1702.

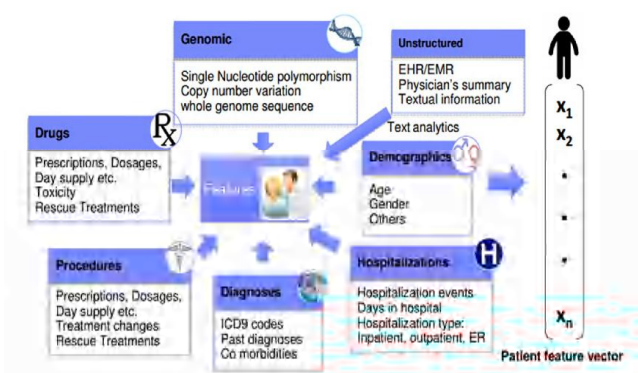
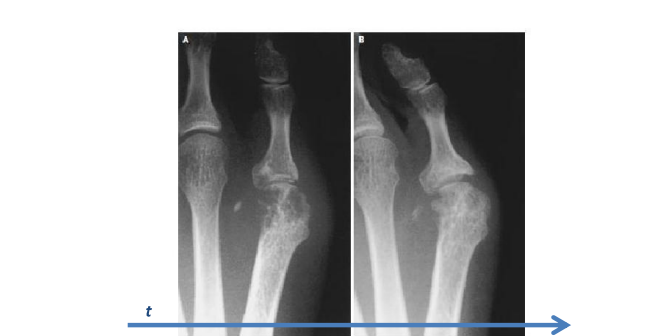


Image credit to Michal Rosen-Zvi

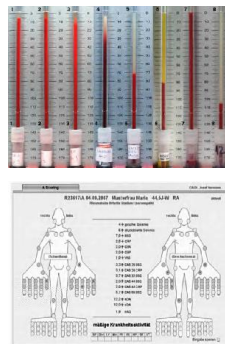


Chao, J., Parker, B. A. & Zvaifler, N. J. (2009) Accelerated Cutaneous Nodulosis Associated with Aromatase Inhibitor Therapy in a Patient with Rheumatoid Arthritis. *The Journal of Rheumatology*, 36, 5, 1087-1088.

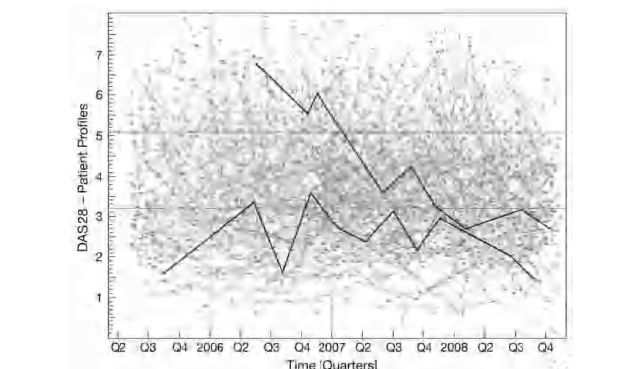


Ikari, K. & Momohara, S. (2005) Bone Changes in Rheumatoid Arthritis. *New England Journal of Medicine*, 353, 15, e13.

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



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



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





<http://biomedicalcomputationreview.org/content/clinical-decision-support-providing-quality-healthcare-help-computer>

Holzinger Group hci-kdd.org

28

Machine Learning Health 08



- Example of a Decision Problem
- Soccer player considering knee surgery
- Uncertainties:
- Success: recovering full mobility
- Risks: infection in surgery (if so, needs another surgery and may lose more mobility)
- Survival chances of surgery

Harvard-MIT Division of Health Sciences and Technology  
HST.951J: Medical Decision Support, Fall 2005

Instructors: Professor Lucila Ohno-Machado and Professor Staal Vinterbo

Holzinger Group hci-kdd.org

31

Machine Learning Health 08

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

The expected utility of decision  $D = d$  is

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

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

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

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



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

Holzinger Group hci-kdd.org

34

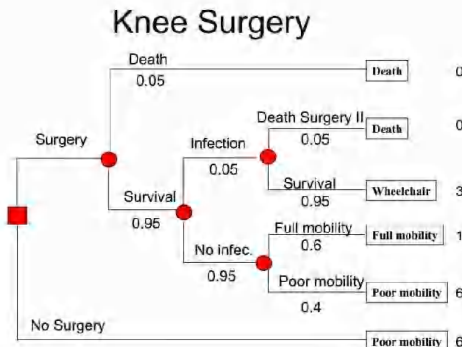
Machine Learning Health 08



Holzinger Group hci-kdd.org

29

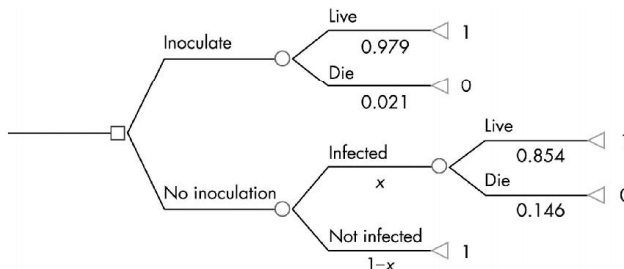
Machine Learning Health 08



Holzinger Group hci-kdd.org

32

Machine Learning Health 08



Ferrando, A., Pagano, E., Scaglione, L., Petrinco, M., Gregori, D. & Ciccone, G. (2009) A decision-tree model to estimate the impact on cost-effectiveness of a venous thromboembolism prophylaxis guideline. *Quality and Safety in Health Care*, 18, 4, 309-313.

Holzinger Group hci-kdd.org

35

Machine Learning Health 08

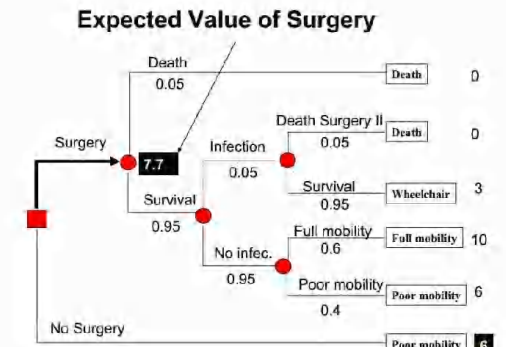
- **Type 1 Decisions:** related to the **diagnosis**, i.e. computers are used to assist in diagnosing a disease on the basis of the individual patient data. Questions include:
  - What is the probability that this patient has a myocardial infarction on the basis of given data (patient history, ECG, ...)?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

Holzinger Group hci-kdd.org

30

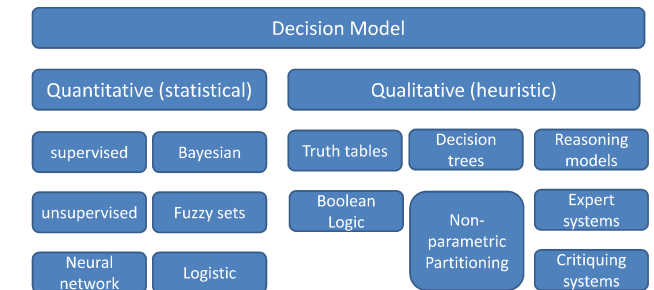
Machine Learning Health 08



Holzinger Group hci-kdd.org

33

Machine Learning Health 08



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

Holzinger Group hci-kdd.org

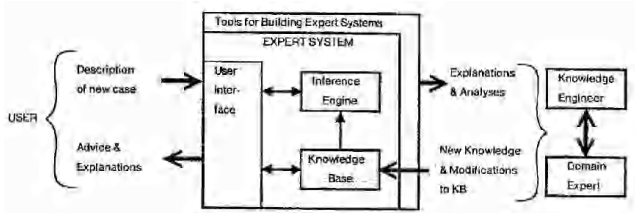
36

Machine Learning Health 08



## 03 History of DSS = History of AI

### Early Knowledge Based System Architecture



Shortliffe, T. & Davis, R. (1975) Some considerations for the implementation of knowledge-based expert systems *ACM SIGART Bulletin*, 55, 9-12.

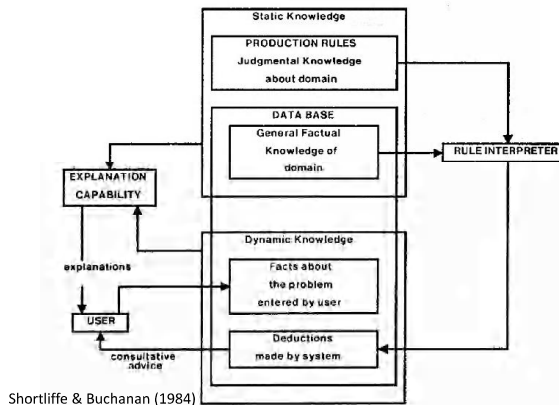
### MYCIN – rule based system - certainty factors

- 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 and 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
- $CF[h] = +1 \rightarrow h$  is 100 % true
- $CF[h] = -1 \rightarrow h$  is 100% false

### A ultrashort history of Early AI

- 1943 McCulloch, W.S. & Pitts, W. A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biology*, 5, (4), 115-133, doi:10.1007/BF02459570.
- 1950 Turing, A.M. Computing machinery and intelligence. *Mind*, 59, (236), 433-460.
- 1959 Samuel, A.L. Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, 3, (3), 210-229, doi:10.1147/rd.33.0210.
- 1975 Shortliffe, E.H. & Buchanan, B.G. 1975. A model of inexact reasoning in medicine. *Mathematical biosciences*, 23, (3-4), 351-379, doi:10.1016/0025-5564(75)90047-4.

### Static Knowledge versus dynamic knowledge



Shortliffe & Buchanan (1984)

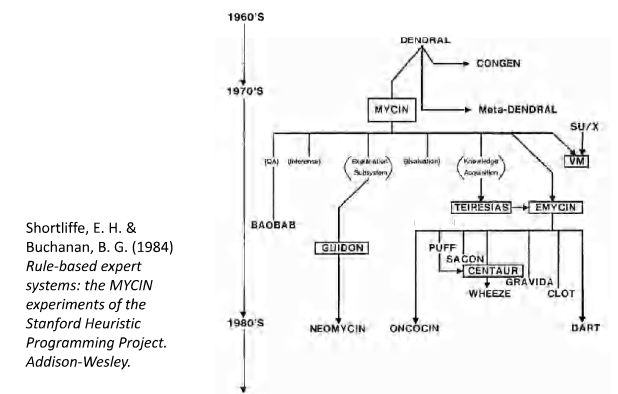
### Original Example from MYCIN

- $h_1$  = 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.

### Evolution of Decision Support Systems (Expert Systems)



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

### 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 was no success in the clinical routine

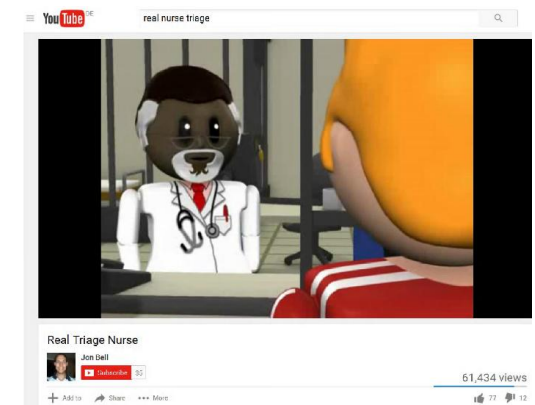




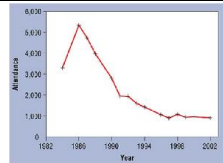
Image credit to Bernhard Schölkopf



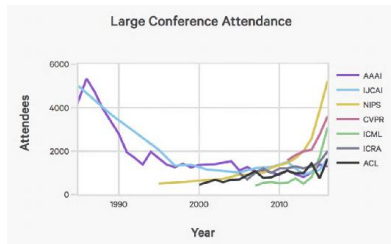
Image credit to Bernhard Schölkopf



<https://blogs.dxc.technology/2017/04/25/are-we-heading-toward-an-ai-winter/>

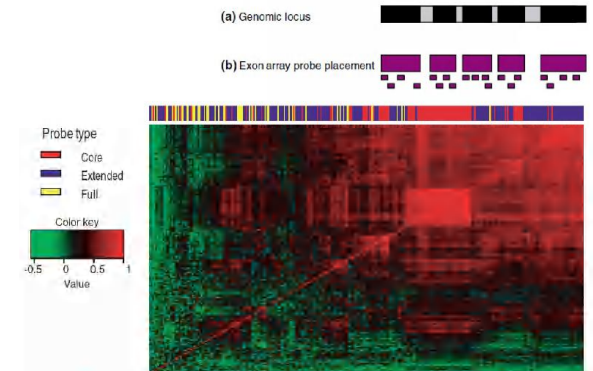


<https://www.computer.org/csli/mags/ex/2003/03/x3018.html>

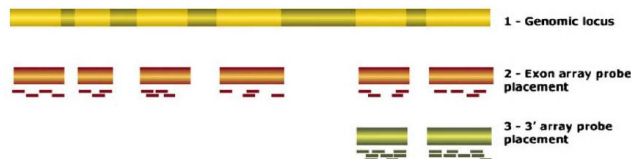


<https://medium.com/machine-learning-in-practice/nips-accepted-papers-stats-26f124843aa0>

## 04 Example: P4-Medicine

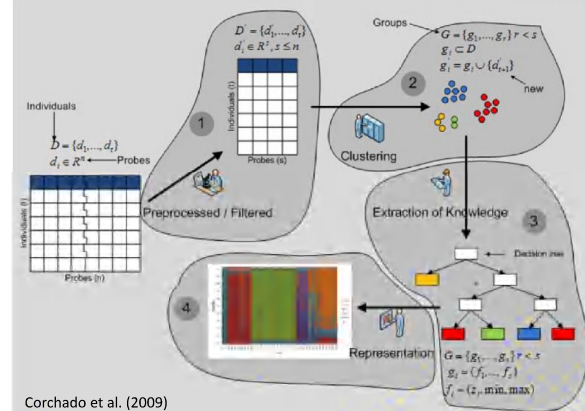


Kapur, K., Xing, Y., Ouyang, Z. & Wong, W. (2007) Exon arrays provide accurate assessments of gene expression. *Genome Biology*, 8, 5, R82.



Exon array structure. Probe design of exon arrays. (1) Exon—intron structure of a gene. Gray boxes represent introns, rest represent exons. Introns are not drawn to scale. (2) Probe design of exon arrays. Four probes target each putative exon. (3) Probe design of 30expression arrays. Probe target the 30end of mRNA sequence.

Corchado, J. M., De Paz, J. F., Rodriguez, S. & Bajo, J. (2009) Model of experts for decision support in the diagnosis of leukemia patients. *Artificial Intelligence in Medicine*, 46, 3, 179-200.

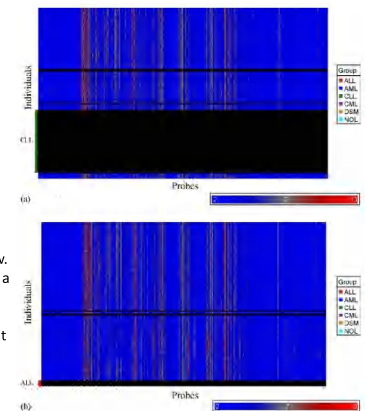


Corchado et al. (2009)

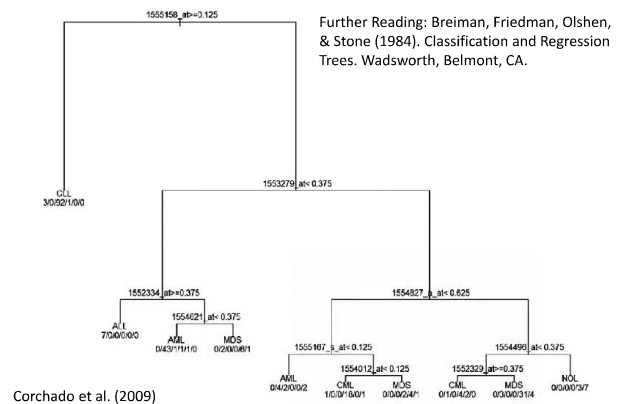
A = acute, C = chronic, L = lymphocytic, M = myeloid

- ALL = cancer of the blood AND bone marrow caused by an abnormal proliferation of lymphocytes.
- AML = cancer in the bone marrow characterized by the proliferation of myeloblasts, red blood cells or abnormal platelets.
- CLL = cancer characterized by a proliferation of lymphocytes in the bone marrow.
- CML = caused by a proliferation of white blood cells in the bone marrow.
- MDS (Myelodysplastic Syndromes) = a group of diseases of the blood and bone marrow in which the bone marrow does not produce a sufficient amount of healthy cells.
- NOL (Normal) = No leukemias

Corchado et al. (2009)

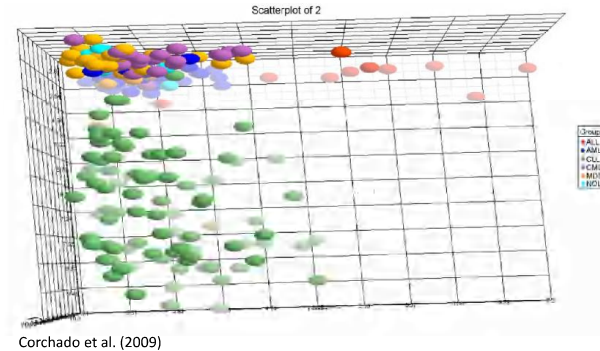






Corchado et al. (2009)

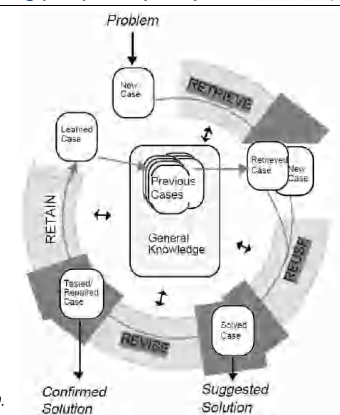
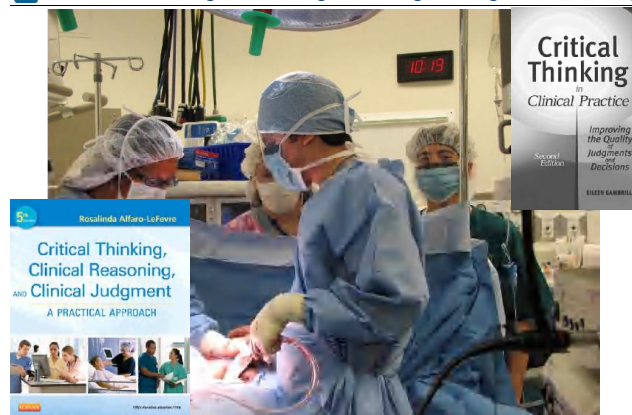
Classification CLL—ALL. Representation of the probes of the decision tree which classify the CLL and ALL to 1555158\_at, 1553279\_at and 1552334\_at



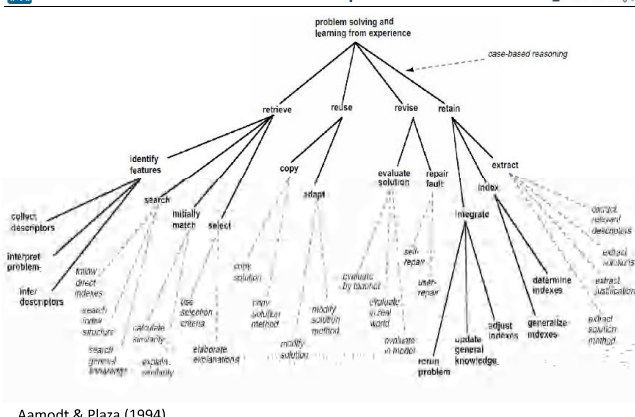
Corchado et al. (2009)

- The model of Corchado et al. (2009) combines:
  - 1) methods to **reduce the dimensionality** of the original data set;
  - 2) pre-processing and data filtering techniques;
  - 3) a clustering method to classify patients; and
  - 4) extraction of knowledge techniques
- The system reflects how human experts work in a lab, but
  - 1) **reduces the time** for making predictions;
  - 2) **reduces the rate of human error**; and
  - 3) **works with high-dimensional data** from exon arrays

## 05 Example: Case Based Reasoning (CBR)



Aamodt, A. & Plaza, E. (1994) Case-based reasoning: Foundational issues, methodological variations, and system approaches. *AI Communications*, 7, 1, 39-59.



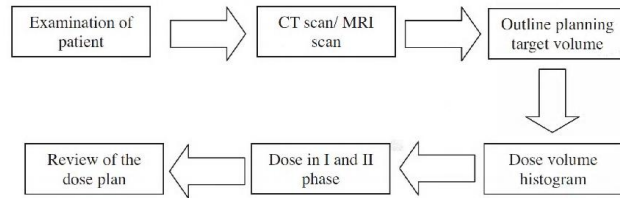
Aamodt & Plaza (1994)



Source: <http://www.teachingmedicalphysics.org.uk>



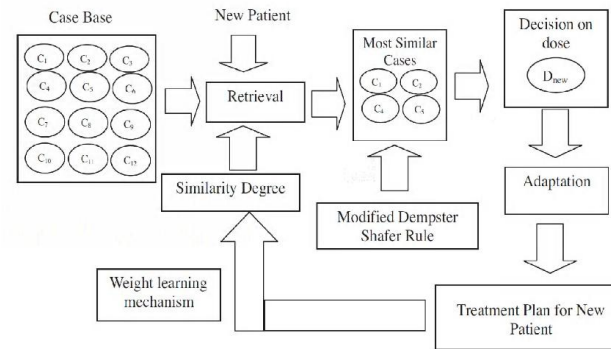
Source: Imaging Performance Assessment of CT Scanners Group, <http://www.impactscan.org>



Measures:

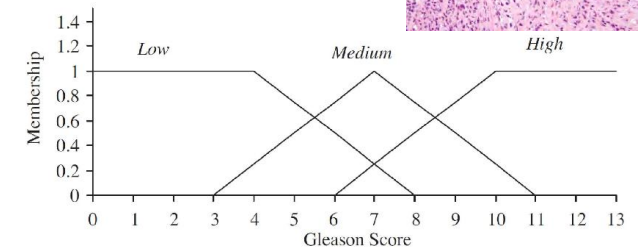
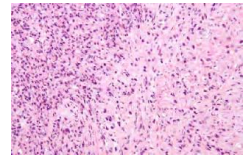
- 1) Clinical Stage = a labelling system
- 2) Gleason Score = grade of prostate cancer = integer between 1 to 10; and
- 3) Prostate Specific Antigen (PSA) value between 1 to 40
- 4) Dose Volume Histogram (DVH) = pot. risk to the rectum (66, 50, 25, 10 %)

Petrovic, S., Mishra, N. & Sundar, S. (2011) A novel case based reasoning approach to radiotherapy planning. *Expert Systems With Applications*, 38, 9, 10759-10769.



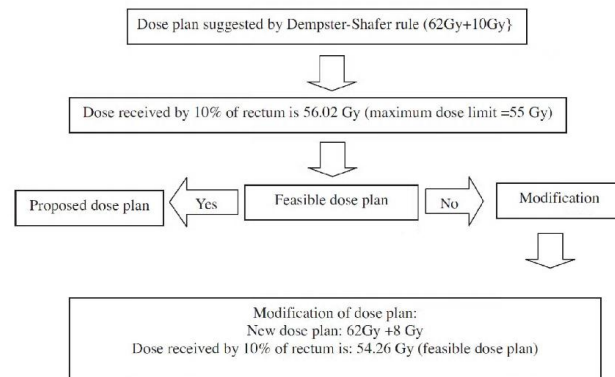
Petrovic, S., Mishra, N. & Sundar, S. (2011) A novel case based reasoning approach to radiotherapy planning. *Expert Systems With Applications*, 38, 9, 10759-10769.

Gleason score evaluates the grade of prostate cancer. Values: integer within the range

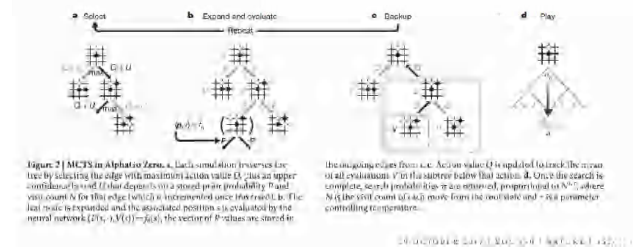


Petrovic, S., Mishra, N. & Sundar, S. (2011) A novel case based reasoning approach to radiotherapy planning. *Expert Systems With Applications*, 38, 9, 10759-10769.

Petrovic et al. (2011)



## 06 Towards Explainable AI



$$(p, v) = f_{\theta}(s) \text{ and } l = (z - v)^2 - \pi^T \log p + c \|\theta\|^2$$

David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, AJ Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George Van Den Driessche, Thore Graepel & Demis Hassabis 2017. Mastering the game of go without human knowledge. *Nature*, 550, (7676), 354-359, doi:10.1038/nature24270.



David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel & Demis Hassabis 2016. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529, (7587), 484-489, doi:10.1038/nature16961.



Andrei Karpathy & Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015. 3128-3137.

Image Captions by deep learning : [github.com/karpathy/neuraltalk2](https://github.com/karpathy/neuraltalk2)

Image Source: Gabriel Villena Fernandez; Agence France-Press, Dave Martin (left to right)

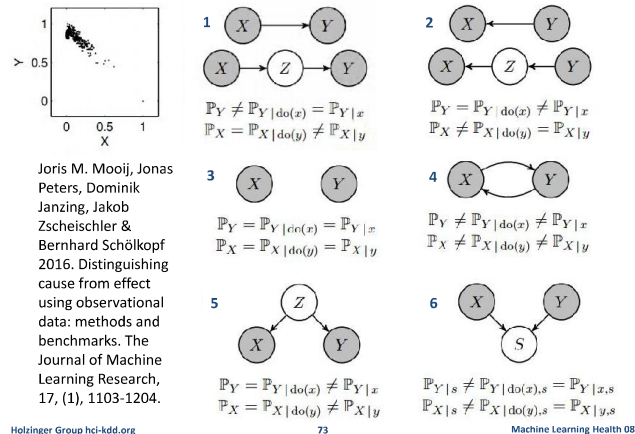
Hans Holbein d.J., 1533, The Ambassadors, London: National Gallery

Lopez-Paz, D., Muandet, K., Schölkopf, B. & Tolstikhin, I. 2015. Towards a learning theory of cause-effect inference. *Proceedings of the 32nd International Conference on Machine Learning*, JMLR, Lille, France.



<https://www.youtube.com/watch?v=9KiVNIUMmCc>

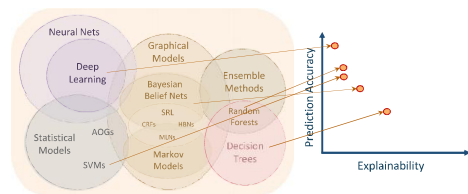
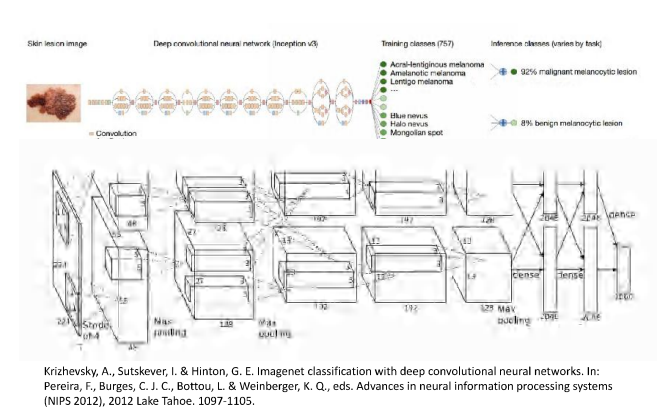




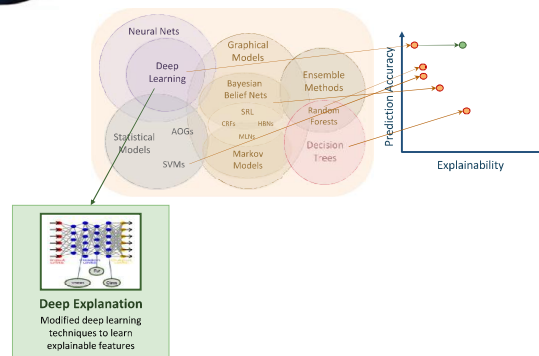
- “How do humans generalize from few examples?”
- Learning relevant representations
- Disentangling the explanatory factors
- Finding the shared underlying explanatory factors, in particular between  $P(x)$  and  $P(Y|X)$ , with a causal link between  $Y \rightarrow X$

Bengio, Y., Courville, A. & Vincent, P. 2013. Representation learning: A review and new perspectives. IEEE transactions on pattern analysis and machine intelligence, 35, (8), 1798-1828, doi:10.1109/TPAMI.2013.50.

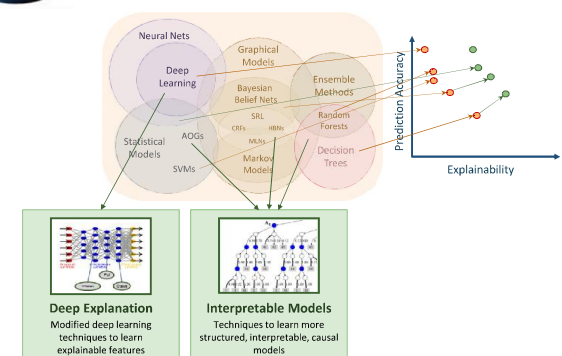
Tenenbaum, J. B., Kemp, C., Griffiths, T. L. & Goodman, N. D. 2011. How to grow a mind: Statistics, structure, and abstraction. Science, 331, (6022), 1279-1285, doi:10.1126/science.1192788.



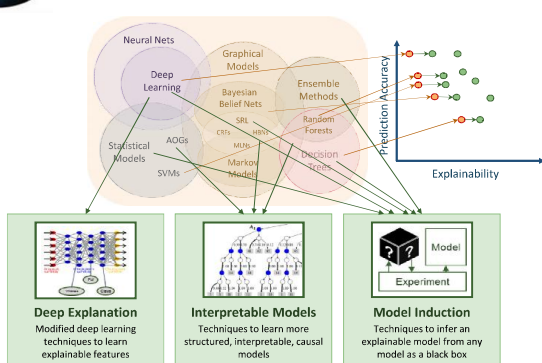
David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



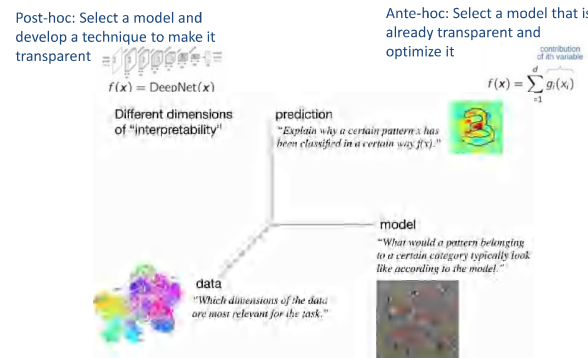
David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



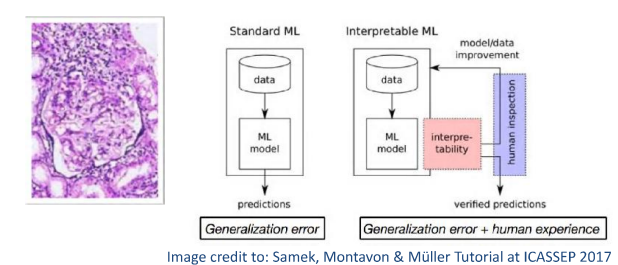
David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



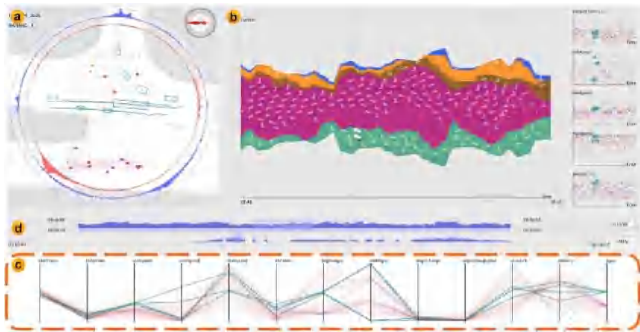
Montavon, G., Samek, W. & Müller, K.-R. 2017. Methods for interpreting and understanding deep neural networks. arXiv:1706.07979.



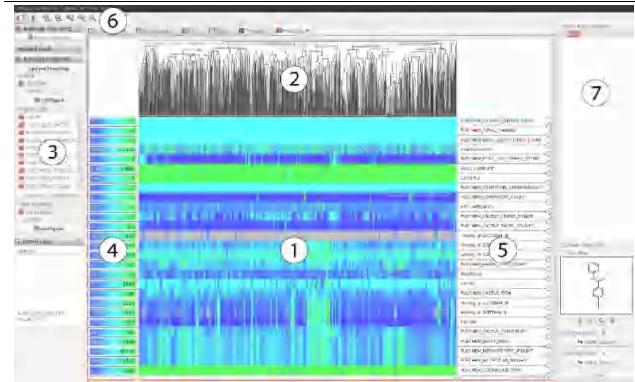
- Wrong decisions can be costly and dangerous!
- Verify that classifier works as expected
- Improve classifier continuously
- Human learning inspired by machine learning

- Interpretability as a novel kind for supporting teaching, learning and knowledge discovery,
- Particularly in abstract fields (informatics)
- Compliance to European Law “the right of explanation”
- Check for bias in machine learning results
- Fostering trust, acceptance, making clear the reliability

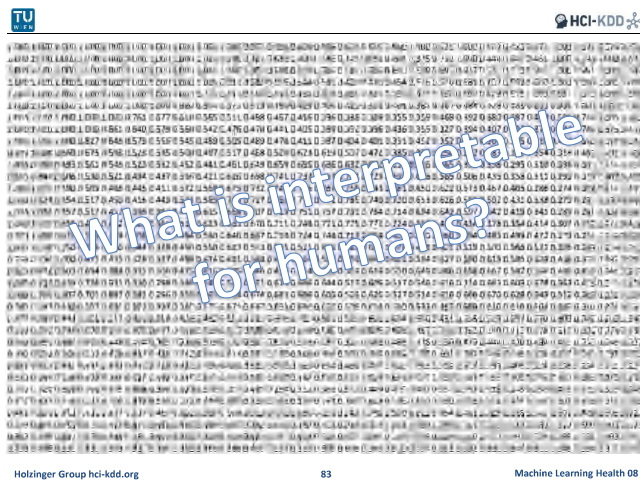
Andreas Holzinger 2018. Explainable AI (ex-AI). Informatik-Spektrum, 41, (2), 138-143, doi:10.1007/s00287-018-1102-5.  
Holzinger Group hci-kdd.org 82 Machine Learning Health 08



Holzinger Group hci-kdd.org 85 Machine Learning Health 08



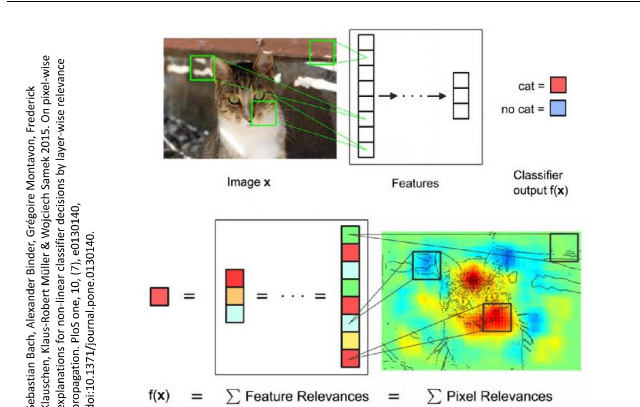
Werner Sturm, Till Schaefer, Tobias Schreck, Andreas Holzinger & Torsten Ullrich. Extending the Scaffold Hunter Visualization Toolkit with Interactive Heatmaps In: Borgo, Rita & Turkey, Cagatay, eds. EG UK Computer Graphics & Visual Computing CGVC 2015, 2015 University College London (UCL). Euro Graphics (EG), 77-84, doi:10.2312/cgvc.20151247.  
Holzinger Group hci-kdd.org 88 Machine Learning Health 08



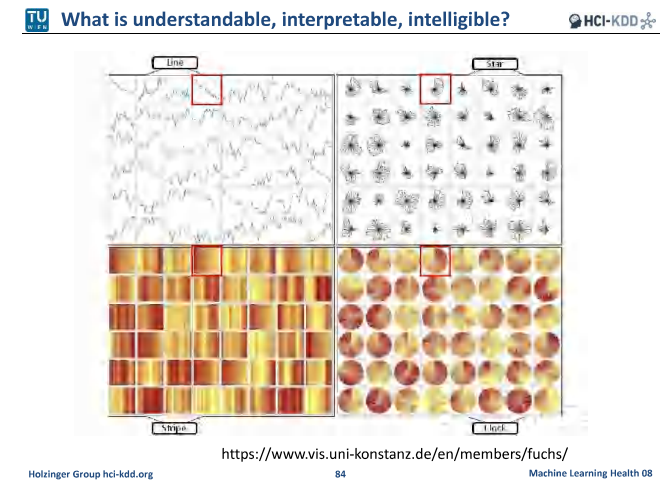
Holzinger Group hci-kdd.org 83 Machine Learning Health 08

## 07 Methods of Explainable AI

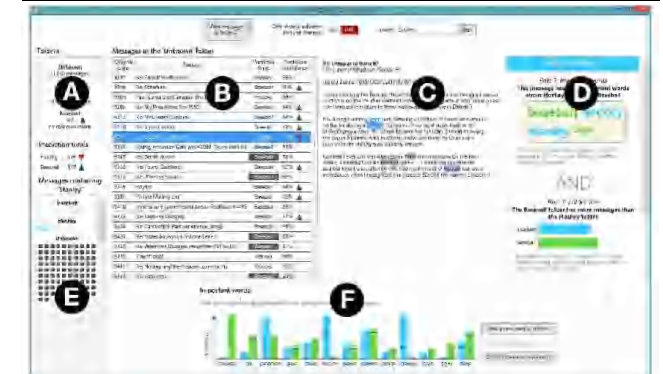
Holzinger Group hci-kdd.org 86 Machine Learning Health 08



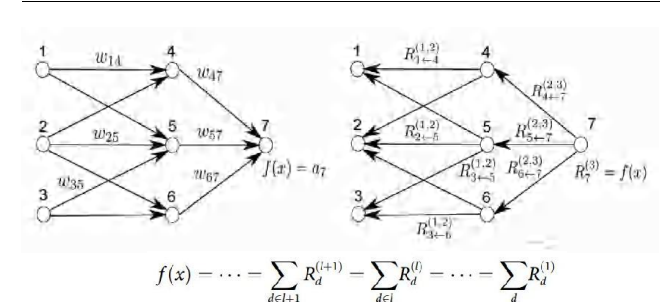
Holzinger Group hci-kdd.org 89 Machine Learning Health 08



Holzinger Group hci-kdd.org 84 Machine Learning Health 08

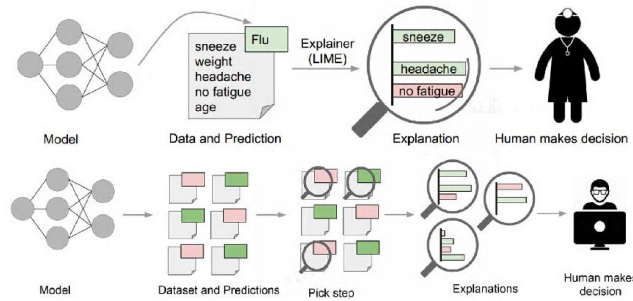


Todd Kulesza, Margaret Burnett, Weng-Keen Wong & Simone Stumpf. Principles of explanatory debugging to personalize interactive machine learning. Proceedings of the 20th International Conference on Intelligent User Interfaces (IUI 2015), 2015 Atlanta. ACM, 126-137, doi:10.1145/2678025.2701399.  
Holzinger Group hci-kdd.org 87 Machine Learning Health 08

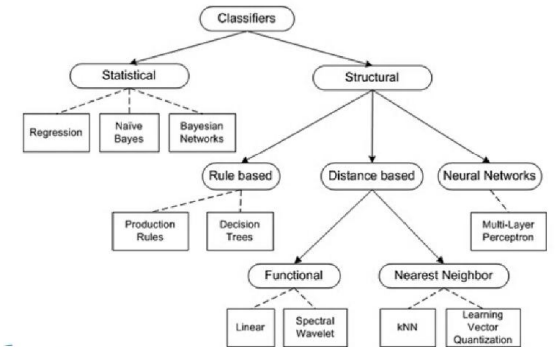


Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.  
Holzinger Group hci-kdd.org 90 Machine Learning Health 08





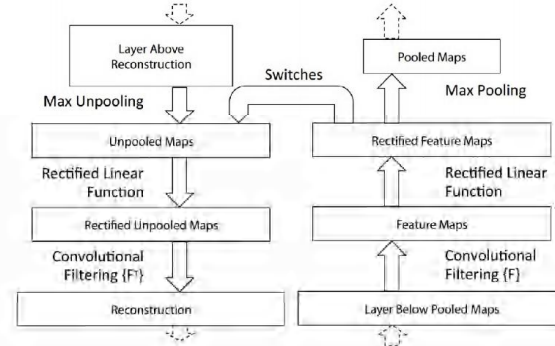
Marco Tulio Ribeiro, Sameer Singh & Carlos Guestrin. Why should I trust you?: Explaining the predictions of any classifier. 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016. ACM, 1135-1144, doi:10.1145/2939672.2939778.



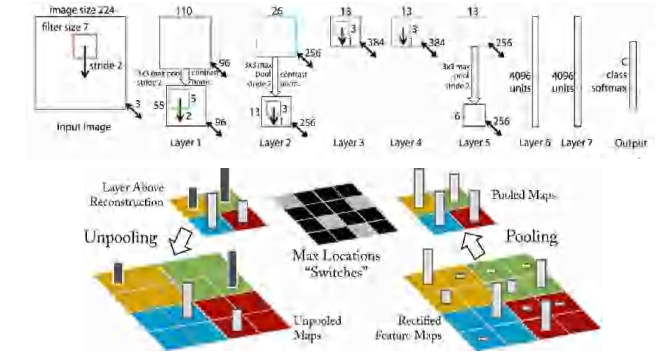
<https://stats.stackexchange.com/questions/271247/machine-learning-statistical-vs-structural-classifiers>



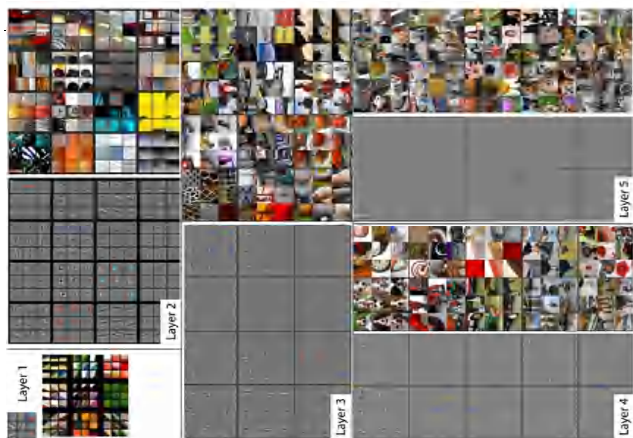
Himabindu Lakkaraju, Ece Kamar, Rich Caruana & Jure Leskovec 2017. Interpretable and Explainable Approximations of Black Box Models. arXiv:1707.01154.



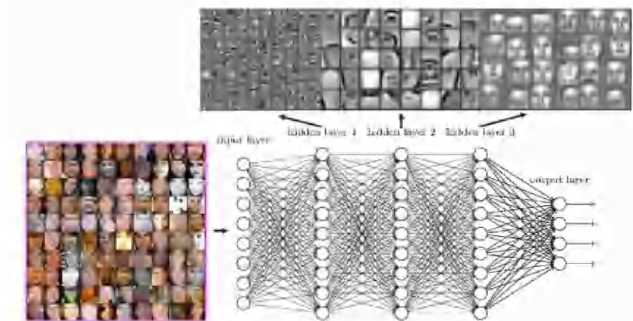
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.



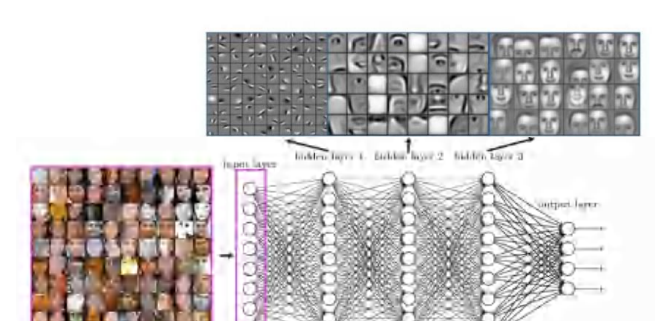
Matthew D. Zeiler & Rob Fergus 2014. Visualizing and understanding convolutional networks. In: D., Fleet, T., Pajdla, B., Schiele & T., Tuytelaars (eds.) ECCV, Lecture Notes in Computer Science LNCS 8689. Cham: Springer, pp. 818-833, doi:10.1007/978-3-319-10590-1\_53.



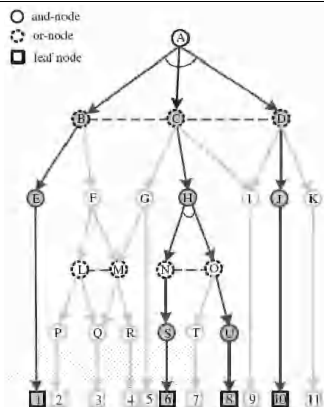
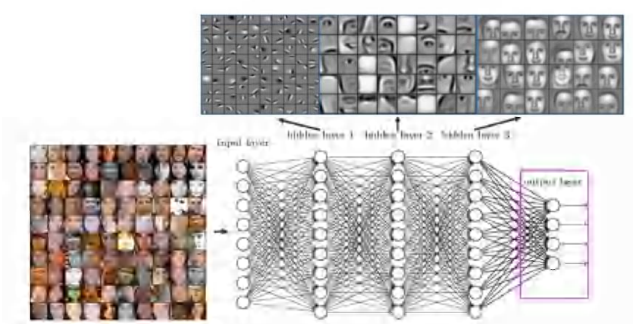
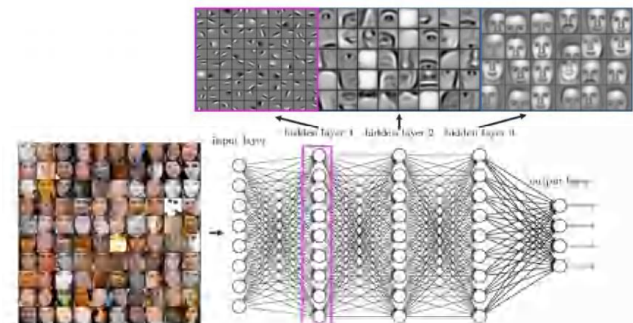
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.



Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901

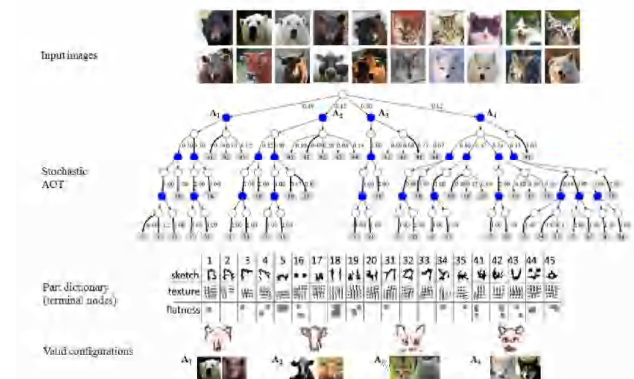
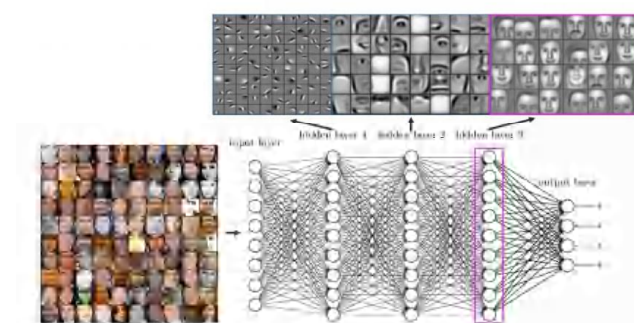
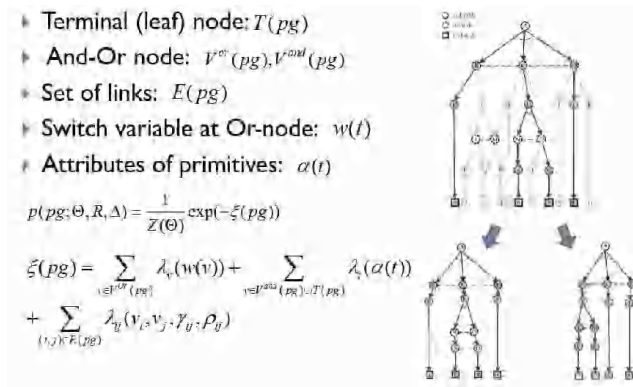
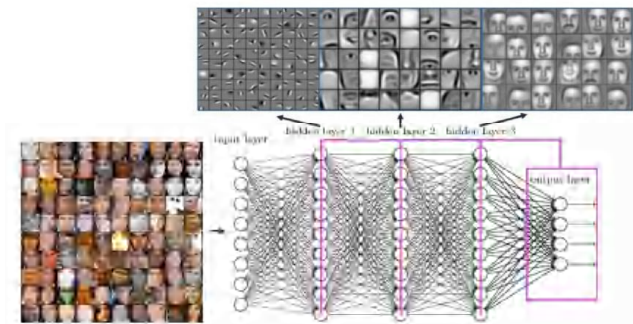
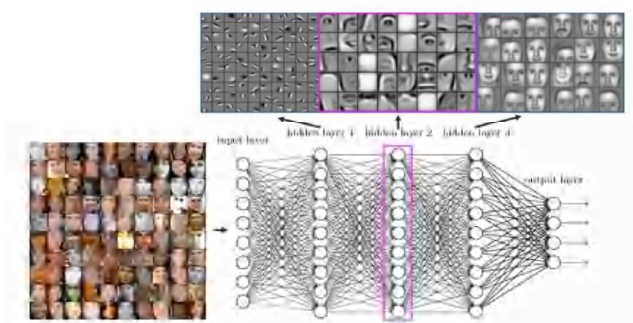




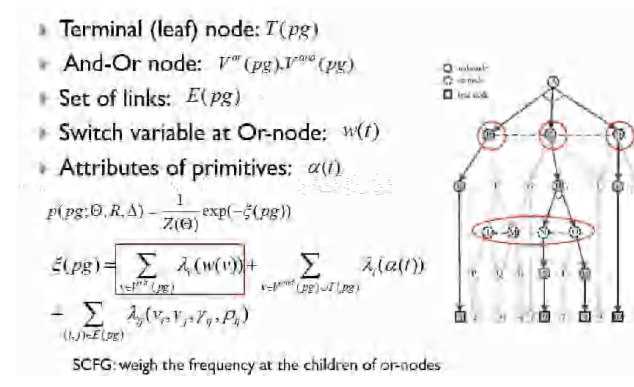


- Algorithm for this framework
  - Top-down/bottom-up computation
- Generalization of small sample
  - Use Monte Carlos simulation to synthesis more configurations
- Fill semantic gap

Images credit to Zhaoyin Jia (2009)



Zhangzhang Si & Song-Chun Zhu 2013. Learning and-or templates for object recognition and detection. IEEE transactions on pattern analysis and machine intelligence, 35, (9), 2189-2205, doi:10.1109/TPAMI.2013.35.





- Terminal (leaf) node:  $T(pg)$
- And-Or node:  $V^{or}(pg), V^{and}(pg)$
- Set of links:  $E(pg)$
- Switch variable at Or-node:  $w(t)$
- Attributes of primitives:  $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\xi(pg) = \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg)} \lambda_v(\alpha(v)) + \sum_{(v,j) \in E(pg)} \lambda_{vj}(v_j, \gamma_{vj}, \rho_{vj})$$

Weigh the local compatibility of primitives (geometric and appearance)

- Terminal (leaf) node:  $T(pg)$
- And-Or node:  $V^{or}(pg), V^{and}(pg)$
- Set of links:  $E(pg)$
- Switch variable at Or-node:  $w(t)$
- Attributes of primitives:  $\alpha(t)$

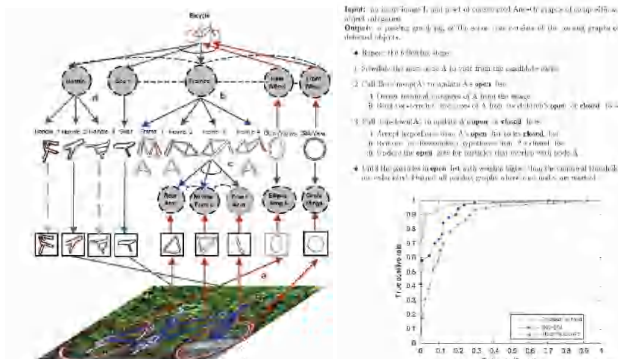
$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\xi(pg) = \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg)} \lambda_v(\alpha(v)) + \sum_{(v,j) \in E(pg)} \lambda_{vj}(v_j, \gamma_{vj}, \rho_{vj})$$

Spatial and appearance between primitives (parts or objects)

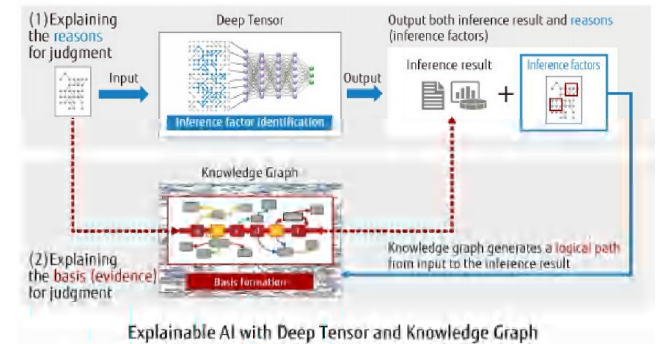
- Terminal (leaf) node:  $T(pg)$
- And-Or node:  $V^{or}(pg), V^{and}(pg)$
- Set of links:  $E(pg)$
- Switch variable at Or-node:  $w(t)$
- Attributes of primitives:  $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\xi(pg) = \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg)} \lambda_v(\alpha(v)) + \sum_{(v,j) \in E(pg)} \lambda_{vj}(v_j, \gamma_{vj}, \rho_{vj})$$


Liang Lin, Tianfu Wu, Jake Porway & Zijian Xu 2009. A stochastic graph grammar for compositional object representation and recognition. Pattern Recognition, 42, (7), 1297-1307, doi:10.1016/j.patcog.2008.10.033.

## Future Work



[http://www.fujitsu.com/jp/Images/artificial-intelligence-en\\_tcm102-3781779.png](http://www.fujitsu.com/jp/Images/artificial-intelligence-en_tcm102-3781779.png)

- What is a good explanation?
- (obviously if the other did understand it)
- Experiments needed!
- What is explainable/understandable/intelligible?
- When is it enough (Sättigungsgrad – you don't need more explanations – enough is enough)
- But how much is it ...



### Teaching Meaningful Explanations

Noel C.F. Codella,\* Michael Hind,\* Karthikeyan Natesan Ramanurthy,\* Murray Campbell, Amit Dhanandhar, Kush R. Varshney, Dennis Wei & Aleksandra Mojsilovic 2018. Teaching Meaningful Explanations. arXiv:1805.11648.

\* These authors contributed equally.

IBM Research  
Yorktown Heights, NY 10598  
(noelcodella, hind, krtavara, krtavara, krtavara, dwei, alexandra@ibm.com)

#### Abstract

The adoption of machine learning in high-stakes applications such as healthcare and law has gained in part because predictions are not accompanied by explanations concerning the underlying domain, which often holds valuable information for decisions and outcomes. In this paper, we propose an approach to generate such explanations in which probing data is represented as features. In addition to features and labels, explanations extract information from domain-specific knowledge to produce both labels and explanations from the joint features. This simple idea ensures that explanations are relevant to the complex explanations and domain knowledge of the consumer. Evaluation spans multiple modeling techniques on a sample game dataset, an image dataset, and a chemical order dataset. Showing that our approach is generalizable across domains and algorithms. Results demonstrate that meaningful explanations can be reliably taught to machine learning algorithms and domain users, together making a step forward.

#### 1 Introduction

These regulations will be amended to require machine systems to provide "meaningful and actionable" information to users in order to ensure transparency. Such and other regulations are part of "meaningful information" to information that should be understandable to the audience consistently. Individuals



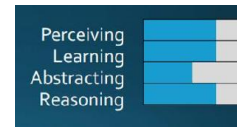
- Engineers create learning models for specific tasks and train them with “big data” (e.g. Deep Learning)
- Advantage: works well for standard classification tasks and has prediction capabilities
- Disadvantage: No contextual capabilities and minimal reasoning abilities

Image credit to John Launchbury



Image credit to John Launchbury

- Computational approaches can find in  $R^n$  what no human is able to see
- However, still there are many hard problems where a human expert in  $R^2$  can understand the **context** and bring in experience, expertise, knowledge, intuition, ...
- Black box approaches can not explain **WHY** a decision has been made ...



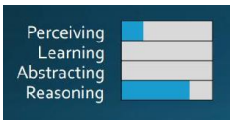
- A contextual model can perceive, learn and understand and abstract and reason
- Advantage: can use transfer learning for adaptation on unknown unknowns
- Disadvantage: Superintelligence ...

Image credit to John Launchbury

- Computers are incredibly fast, accurate and stupid,
- humans are incredibly slow, inaccurate and brilliant,
- together they are powerful beyond imagination

(Einstein never said that)

<https://www.benshoemate.com/2008/11/30/einstein-never-said-that>



- Engineers create a set of logical rules to represent knowledge (Rule based Expert Systems)
- Advantage: works well in narrowly defined problems of well-defined domains
- Disadvantage: No adaptive learning behaviour and poor handling of  $p(x)$

Image credit to John Launchbury

- Myth 1a: Superintelligence by 2100 is inevitable!
- Myth 1b: Superintelligence by 2100 is impossible!
- Fact: We simply don't know it!**
- Myth 2: Robots are our main concern  
**Fact: Cyberthreats are the main concern: it needs no body – only an Internet connection**
- Myth 3: AI can never control us humans  
**Fact: Intelligence is an enabler for control: We control tigers by being smarter ...**



Thank you!