



Assoc.Prof. Dr. Andreas Holzinger  
185.A83 Machine Learning for Health Informatics  
2019S, VU, 2.0 h, 3.0 ECTS  
Lecture 01 - Dienstag, 12.03.2019



## MAKE Health Machine Learning & Knowledge Extraction in health informatics: challenges & directions

andreas.holzinger AT tuwien.ac.at  
<https://hci-kdd.org/machine-learning-for-health-informatics-class-2019>



Holzinger Group hci-kdd.org

1

2019 Machine Learning for Health 01

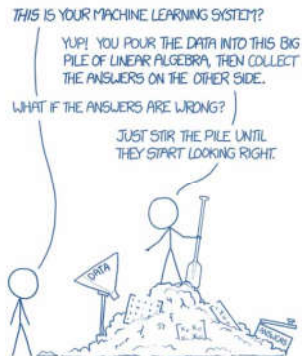


Image Source: Randall Munroe <https://xkcd.com>

Holzinger Group hci-kdd.org

4

2019 Machine Learning for Health 01

- algorithm development is at the core – however, successful ML needs a concerted effort of various topics ...



Holzinger Group hci-kdd.org

7

2019 Machine Learning for Health 01

## LV 185.A83 Machine Learning for Health Informatics (Class of 2019)

Study Code: 066 936 Master program Medical Informatics

<https://iss.tuwien.ac.at/curriculum/public/curriculum.xhtml?dowid=9468&doid=253&key=56089&semester=NEXT>

Semester hours: 2.0 h; ECTS-Credits: 3.0; Type: VU Lecture and Exercise

ECTS-Breakdown (sum=75 h, corresponds with 3 ECTS, where 1 ECTS = 25 h workload):

Presence during lecture	8 * 3 h	24 h
Preparation before and after lecture	8 * 1 h	08 h
Preparation of assignments and presentation	28 h + 2	30 h
Written exam including preparation	1 h + 12 h	13 h
TOTAL students' workload		75 h

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

All Slides will be put on-line AFTER each class!

Holzinger Group hci-kdd.org

2

2019 Machine Learning for Health 01

- 01 The HCI-KDD approach: integrative ML
- 02 Application Area Health
- 03 Probabilistic Learning
- 04 Automatic Machine Learning (aML)
- 05 Interactive Machine Learning (iML)
- 06 Causality vs. Causability
- 07 Explainable AI
- Conclusion and Future Outlook

Holzinger Group hci-kdd.org

5

2019 Machine Learning for Health 01

Class Schedule for 2019 (subject to change: please check class URL for any changes):

Nr	Day, Date	Time	h	Topic
1	Dienstag 12.3.2019	17:30- 20:30	3 h	Machine learning for health informatics: Introduction, challenges and future directions
2	Dienstag 19.3.2019	17:30- 20:30	3h	From clinical decision making to explainable AI: selected methods of transparent machine learning
3	Dienstag 26.3.2019	17:30- 20:30	3 h	Tutorial: Augmentation and Explainability And FIRST ASSIGNMENT
4	Dienstag 02.4.2019	17:30- 20:30	3 h	Probabilistic Graphical Models: from knowledge representation to graph model learning
5	Dienstag 09.4.2019	17:30- 20:30	3 h	Tutorial: Probabilistic Programming with Python and SECOND ASSIGNMENT
Easter Break and Time for working on the assignments				
6	Dienstag 30.4.2019	17:30- 20:30	3 h	Data for machine learning: quality, fusion, integration, probabilistic information and entropy
7	Dienstag 07.5.2019	17:30- 20:30	3 h	Causality and causal machine learning for decision support, ethical, legal and social issues of AI in health
Finalization of assignments				
8	Dienstag 28.5.2019	17:30- 20:30	3 h	Final exam (written test, 40 %) and presentations of the assignments (orally, 30 %) quality of the assignments 25 % each (coding, 50 %)

Transparent Procedure how to  
get grades: sample exam will be  
made openly available

Holzinger Group hci-kdd.org

3

2019 Machine Learning for Health 01

## 01 What is the HCI-KDD approach?

Holzinger Group hci-kdd.org

6

2019 Machine Learning for Health 01

- algorithm development is at the core – however, successful ML needs a concerted effort of various topics ...



Andreas Holzinger 2017. Introduction to Machine Learning and Knowledge Extraction (MAKE). *Machine Learning and Knowledge Extraction*, 1, (1), 1-20, doi:10.3390/make1010001.

Holzinger Group hci-kdd.org

8

2019 Machine Learning for Health 01



<http://www.bach-cantatas.com>

Holzinger Group hci-kdd.org

9

2019 Machine Learning for Health 01



## “Solve intelligence – then solve everything else”

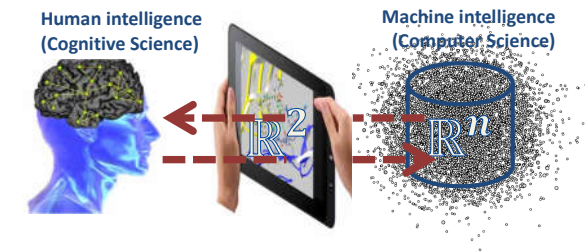


<https://youtu.be/XAbLn66iHcQ?t=1h28m54s>



- 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

## Our goal: Understanding Context !



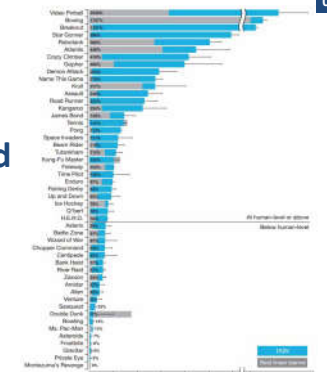
Andreas Holzinger 2013. Human-Computer Interaction and Knowledge Discovery (HCI-KDD): What is the benefit of bringing those two fields to work together? In: Lecture Notes in Computer Science LNCS 8127, pp. 319-328, doi:10.1007/978-3-642-40511-2\_22.



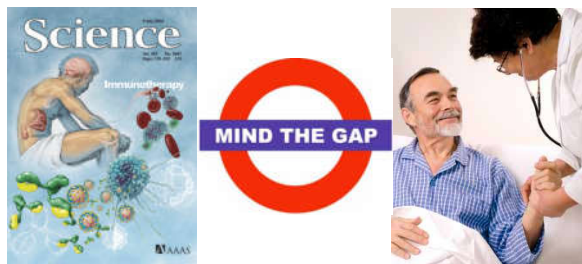
## Why is this application area complex ?

## Compare your best ML algorithm with a seven year old child ...

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







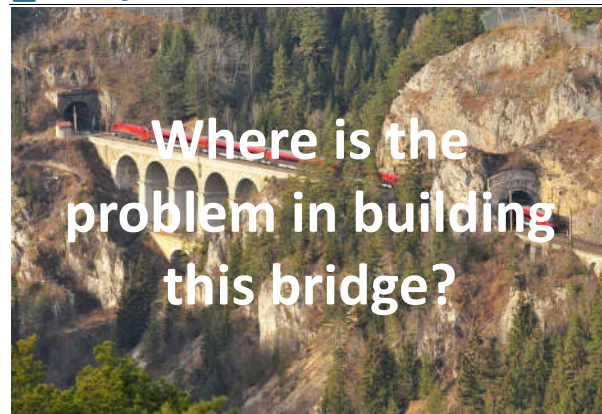
## Our central hypothesis: Information may bridge this gap

Andreas Holzinger & Klaus-Martin Simon (eds.) 2011. Information Quality in e-Health. Lecture Notes in Computer Science LNCS 7058, Heidelberg, Berlin, New York: Springer, doi:10.1007/978-3-642-25364-5.

# 03 Probabilistic Learning



- Newton, Leibniz, ... developed calculus – mathematical language for describing and dealing with rates of change
- Bayes, Laplace, ... developed probability theory - the mathematical language for describing and dealing with uncertainty
- Gauss generalized those ideas



Probability theory is nothing but common sense reduced to calculation ...

$$\hat{y} = \hat{f}(x) = \arg\max_{c=1}^C p(y = c | x, \mathcal{D})$$



Pierre Simon de Laplace (1749-1827)

$$p(x_i) = \sum P(x_i, y_j)$$

$$p(x_i, y_j) = p(y_j | x_i) P(x_i)$$

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

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

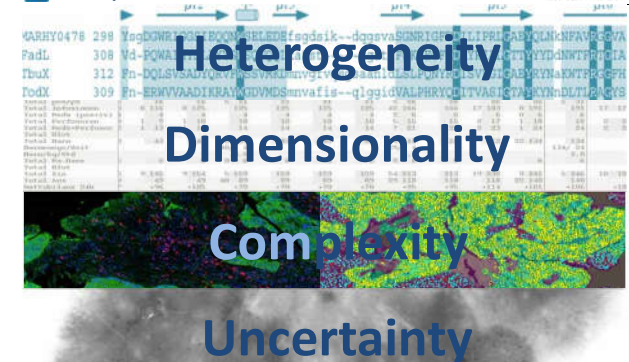
$$P(\text{hypothesis} | \text{data}) = \frac{P(\text{hypothesis}) P(\text{data} | \text{hypothesis})}{\sum_h P(h) P(\text{data} | h)} \quad P(\theta | \mathcal{D}, m) = \frac{P(\mathcal{D} | \theta, m) P(\theta | m)}{P(\mathcal{D} | m)}$$

$P(\mathcal{D} | \theta, m)$  likelihood of parameters  $\theta$  in model  $m$

$P(\theta | m)$  prior probability of  $\theta$

$P(\theta | \mathcal{D}, m)$  posterior of  $\theta$  given data  $\mathcal{D}$

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



Holzinger, A., Dehmer, M. & Jurisica, I. 2014. Knowledge Discovery and interactive Data Mining in Bioinformatics - State-of-the-Art, future challenges and research directions. *BMC Bioinformatics*, 15, (S6), 11. Holzinger Group hci-kdd.org

What is the simplest mathematical operation for us?

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

How do we call repeated adding?

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

Laplace (1773) showed that we can write:

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

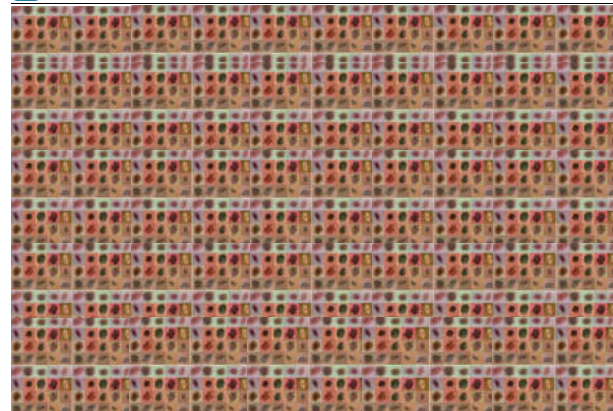
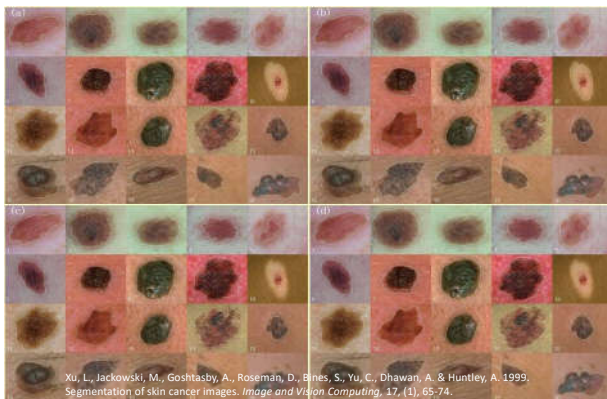
Now we introduce a third, more complicated operation:


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

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

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





$$\mathcal{D} = x_{1:n} = \{x_1, x_2, \dots, x_n\} \quad p(\mathcal{D}|\theta)$$


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

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

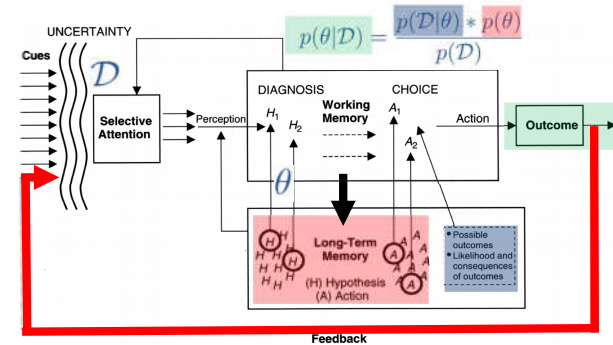
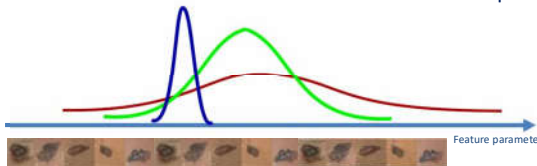
**The inverse probability allows to learn from data, infer unknowns, and make predictions**

$$\begin{array}{lll} d \dots \textit{data} & \mathcal{H} \dots \{H_1, H_2, \dots, H_n\} & \forall h, d \dots \\ h \dots \textit{hypotheses} & & \end{array}$$

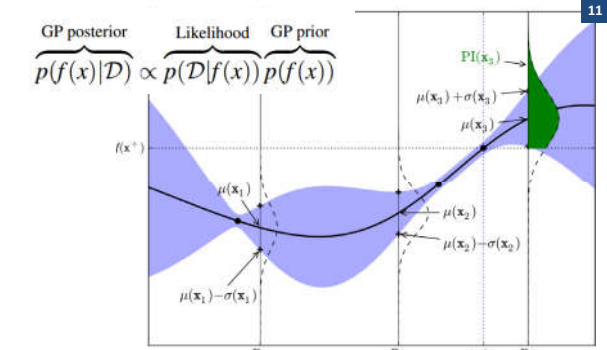
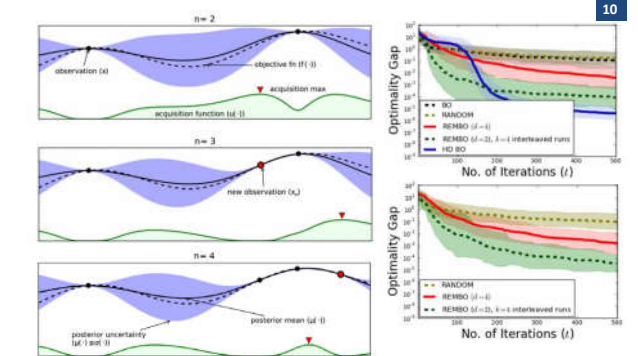
$$p(h|d) = \frac{p(d|h) * p(h)}{\sum_{h \in H} p(d|h') p(h')}$$

## Posterior Probability

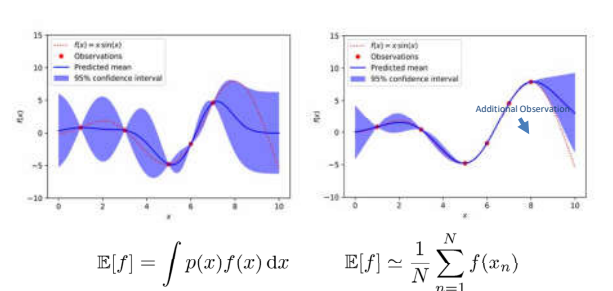
### Problem in $\mathbb{R}^n \rightarrow \text{complex}$



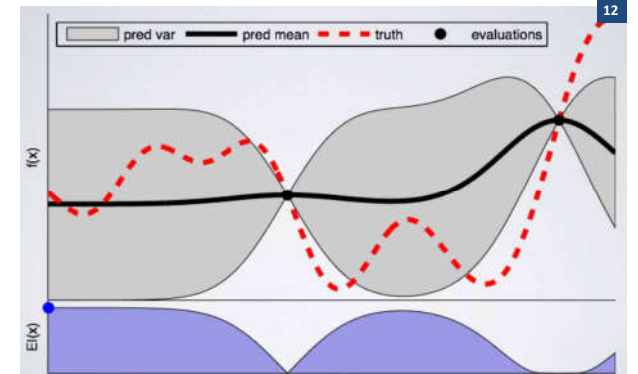
Wickens, C. D. (1984) *Engineering psychology and human performance*. Columbus (OH), Charles Merrill, modified by Holzinger, A.



Brochu, E., Cora, V. M. & De Freitas, N. 2010. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. arXiv:1012.2599.

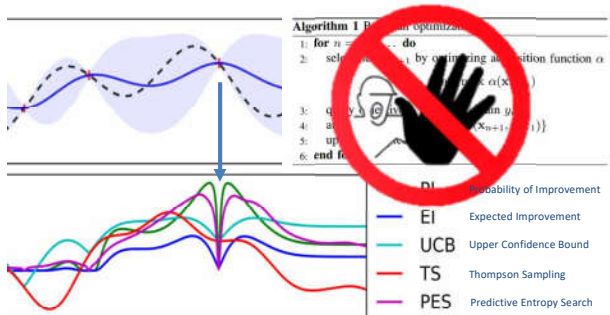
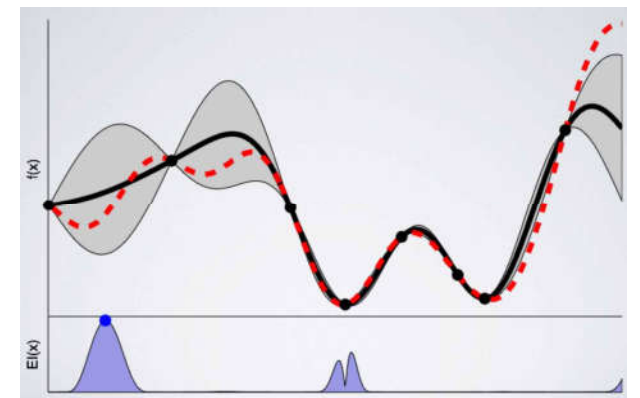
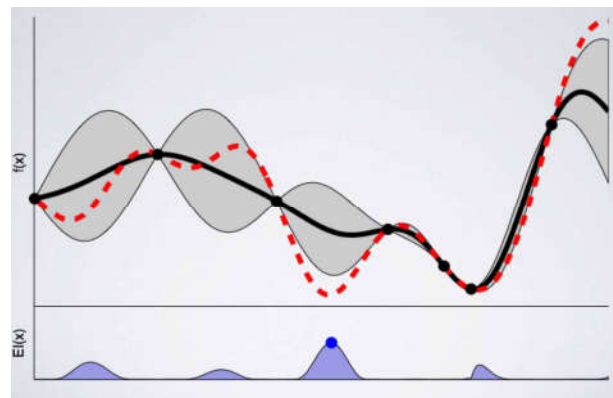
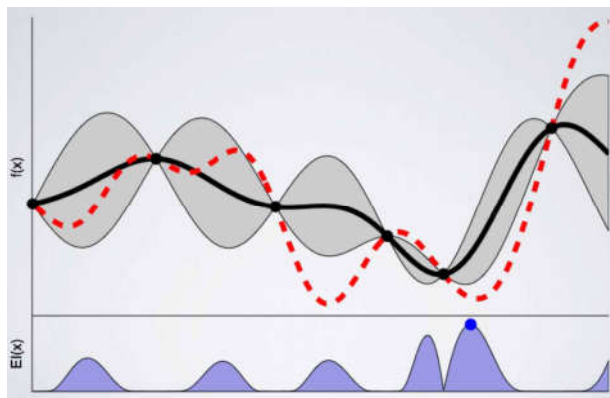
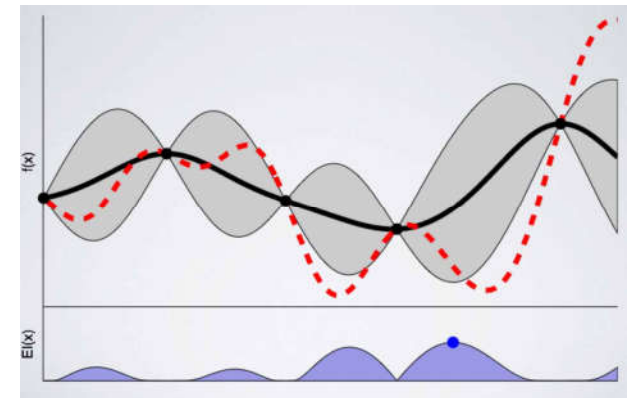
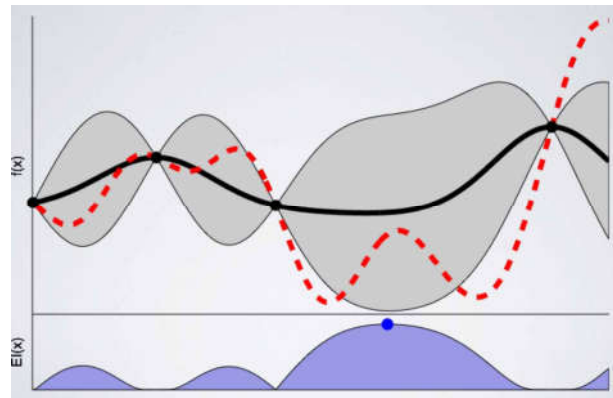
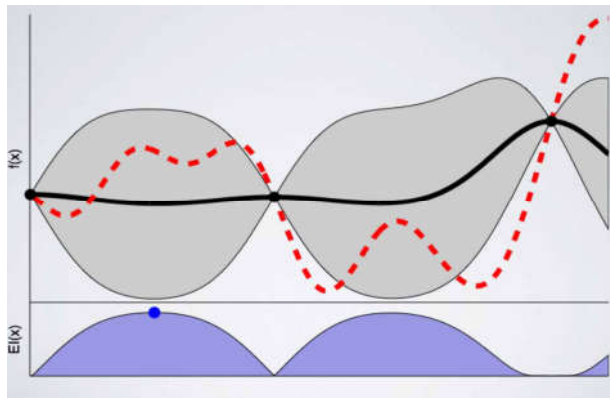


Holzinger, A. 2017. Introduction to Machine Learning and Knowledge Extraction (MAKE). Machine Learning and Knowledge Extraction, 1, (1), 1-20, doi:10.3390/make1010001.



Snoek, J., Larochelle, H. & Adams, R. P. Practical Bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 2012. 2951-2959.





04 aML

Best practice examples of aML ...

**Recommender Systems**

Holzinger Group hci-kdd.org 46 2019 Machine Learning for Health 01

**Fully automatic autonomous vehicles ("Google car")**

Human			Machine		
LEVEL 0	LEVEL 1	LEVEL 2	LEVEL 3	LEVEL 4	LEVEL 5
No Active Assistance System	Longitudinal or Transverse Guide	Longitudinal and Transverse Guide	Assessments for Take Over	No Driver Intervention	No Driver
Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	No Take Over Request	No Take Over Request
Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	No Take Over Request	No Take Over Request
Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	No Take Over Request	No Take Over Request
Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	No Take Over Request	No Take Over Request
Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	Requires full Eyes On	No Take Over Request	No Take Over Request

Guizzo, E. 2011. How google's self-driving car works. IEEE Spectrum Online, 10, 18.

Holzinger Group hci-kdd.org 47 2019 Machine Learning for Health 01

**Why did the car do that? Who is responsible?**

Image Source: <http://www.businessinsider.de/who-is-responsible-when-a-driverless-car-crashes-2016-2?r=US&IR=T>

Holzinger Group hci-kdd.org 48 2019 Machine Learning for Health 01

**... and thousands of industrial aML applications ...**

Seshia, S. A., Juniwal, G., Sadigh, D., Donze, A., Li, W., Jensen, J. C., Jin, X., Deshmukh, J., Lee, E. & Sastry, S. 2015. Verification by, for, and of Humans: Formal Methods for Cyber-Physical Systems and Beyond. Illinois ECE Colloquium.

Holzinger Group hci-kdd.org 49 2019 Machine Learning for Health 01

**Big Data is necessary for aML !**

Sonnenburg, S., Rätsch, G., Schäfer, C. & Schölkopf, B. 2006. Large scale multiple kernel learning. Journal of Machine Learning Research, 7, (7), 1531-1565.

Holzinger Group hci-kdd.org 50 2019 Machine Learning for Health 01

**10 million 200 x 200 px images downloaded from Web**

$$x^* = \arg \min_x f(x; W, H), \text{ subject to } \|x\|_2 = 1.$$

Le, Q. V., Ranzato, M. A., Monga, R., Devin, M., Chen, K., Corrado, G. S., Dean, J. & Ng, A. Y. 2011. Building high-level features using large scale unsupervised learning. arXiv preprint arXiv:1112.6209.

Le, Q. V. 2013. Building high-level features using large scale unsupervised learning. *IEEE Intl. Conference on Acoustics, Speech and Signal Processing ICASSP*. IEEE. 8595-8598, doi:10.1109/ICASSP.2013.6639343.

Holzinger Group hci-kdd.org 51 2019 Machine Learning for Health 01

**Deep Convolutional Neural Network Pipeline**

Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C. J. C., Bottou, L. & Weinberger, K. Q., eds. *Advances in neural information processing systems* (NIPS 2012), 2012 Lake Tahoe. 1097-1105.

Holzinger Group hci-kdd.org 52 2019 Machine Learning for Health 01

**When does aML fail ...**

- Sometimes we do not have "big data", where aML-algorithms benefit.
- Sometimes we have
  - Small amount of data sets
  - Rare Events – no training samples
  - NP-hard problems, e.g.
    - Subspace Clustering,
    - k-Anonymization,
    - Protein-Folding, ...

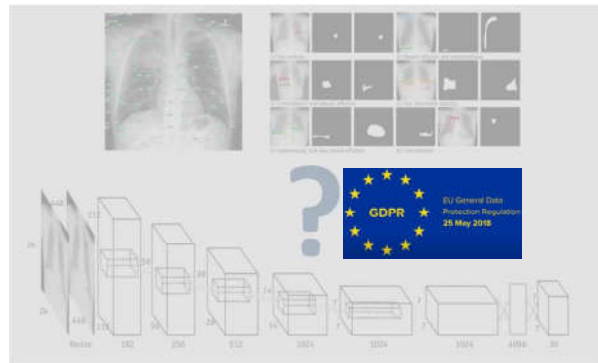
Holzinger Group hci-kdd.org 53 2019 Machine Learning for Health 01

**Houston, we have a problem ...**

Source: NASA, Image is in the public domain

Holzinger Group hci-kdd.org 54 2019 Machine Learning for Health 01





June-Goo Lee, Sanghoon Jun, Young-Won Cho, Hyunna Lee, Guk Bae Kim, Joon Beom Seo & Namkug Kim 2017. Deep learning in medical imaging: general overview. Korean journal of radiology, 18, (4), 570-584, doi:10.3348/kjr.2017.18.4.570.

Holzinger Group hci-kdd.org

55

2019 Machine Learning for Health 01

## There is an urgent need for “explainability”

Holzinger Group hci-kdd.org

56

2019 Machine Learning for Health 01

## 05 iML

Holzinger Group hci-kdd.org

57

2019 Machine Learning for Health 01

16

- iML := algorithms which interact with agents\*) and can optimize their learning behaviour through this interaction
- \*) where the agents can be human

Holzinger, A. 2016. Interactive Machine Learning (iML). Informatik Spektrum, 39, (1), 64-68, doi:10.1007/s00287-015-0941-6.

Holzinger Group hci-kdd.org

58

2019 Machine Learning for Health 01



Image Source: 10 Ways Technology is Changing Healthcare <http://newhealthpost.com> Posted online on April 22, 2018

Holzinger Group hci-kdd.org

59

2019 Machine Learning for Health 01



Holzinger Group hci-kdd.org

60

2019 Machine Learning for Health 01



Holzinger Group hci-kdd.org

61

2019 Machine Learning for Health 01

## Why using human intuition?

Holzinger Group hci-kdd.org

62

2019 Machine Learning for Health 01

### Humans can generalize even from few examples ...

- They learn relevant representations
- Can disentangle the explanatory factors
- Find 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.

Holzinger Group hci-kdd.org

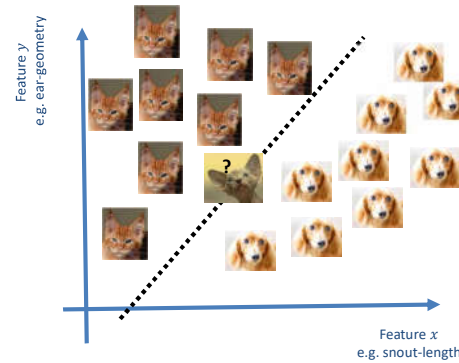
63

2019 Machine Learning for Health 01

## ... can infer from little data ...



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.



See also: Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy 2014. Explaining and harnessing adversarial examples. arXiv:1412.6572.

### Adversarial Examples that Fool both Computer Vision and Time-Limited Humans

Gamaleldin F. Elsayed\*  
Google Brain  
gamaleldin.elsayed@gmail.com

Shreya Shankar  
Stanford University

Brian Cheung  
UC Berkeley

Nicolas Papernot  
Pennsylvania State University

Alex Kurakin  
Google Brain

Ian Goodfellow  
Google Brain

Jascha Soth-Dickstein  
Google Brain  
jaschaad@google.com

#### Abstract

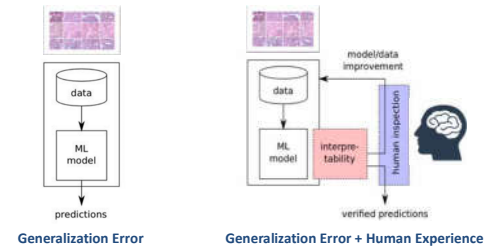
Machine learning models are vulnerable to **adversarial examples**: small changes to images can cause computer vision models to make mistakes such as identifying a school bus as an ostrich. However, it is still an open question whether humans are prone to similar mistakes. Here, we address this question by leveraging recent techniques that transfer adversarial examples from computer vision models with known parameters and architecture to other models with unknown parameters and architecture, and by matching the initial processing of the human visual system. We find that adversarial examples that strongly transfer across computer vision models influence the classifications made by time-limited human observers.

Gamaleldin F Elsayed, Shreya Shankar, Brian Cheung, Nicolas Papernot, Alex Kurakin, Ian Goodfellow & Jascha Soth-Dickstein 2018. Adversarial Examples that Fool both Human and Computer Vision. arXiv:1802.08195.

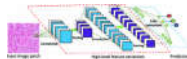


- From black-box to glass-box ML
- Exploit human intelligence for solving hard problems (e.g. Subspace Clustering, k-Anonymization, Protein-Design)
- Towards multi-agent systems with humans-in-the-loop

Holzinger, A., Plass, M., Holzinger, K., Crisan, G., Pintea, C. & Palade, V. 2016. Towards interactive Machine Learning (iML): Applying Ant Colony Algorithms to solve the Traveling Salesman Problem with the Human-in-the-Loop approach. Springer Lecture Notes in Computer Science LNCS 9817. Heidelberg, Berlin, New York: Springer, pp. 81-95, doi:10.1007/978-3-319-45507-56.



Verify that algorithms/classifiers work as expected  
Wrong decisions can be costly and dangerous



Understanding the weaknesses and errors of the ML-Model - Detection of bias in both directions



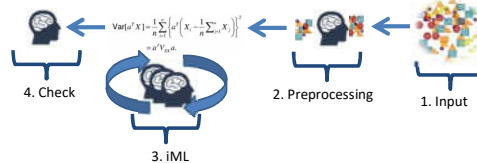
Scientific interpretability, replicability, causality  
The "why" is often more important than the prediction



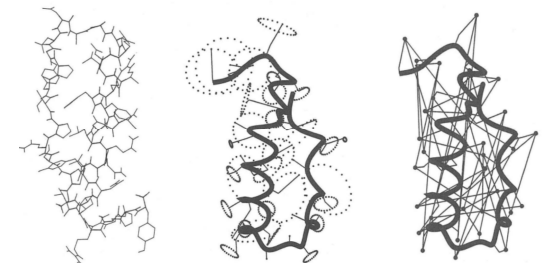
Enable re-traceability, re-enactivity  
Compliance to legislation "right for explanation", retain human reliability, fosters trust and acceptance



**Interactive Machine Learning:** Human is seen as an agent involved in the actual learning phase, step-by-step influencing measures such as distance, cost functions ...



Holzinger, A. 2016. Interactive Machine Learning for Health Informatics: When do we need the human-in-the-loop? Brain Informatics (BRIN), 3, (2), 119-131, doi:10.1007/s40708-016-0042-6.



Bohr, H. & Brunak, S. 1989. A travelling salesman approach to protein conformation. Complex Systems, 3, 9-28



```

Input : ProblemSize, m, β, ρ, σ, q0
Output: Pbest
Pbest ← CreateHeuristicSolution(ProblemSize);
Pbestcost ← Cost(Pbest);
Pheromoneinit ←  $\frac{1}{\beta}$ ;
Pheromone ← InitializePheromone(Pheromoneinit);
while ¬StopCondition() do
  for i = 1 to m do
    Si ← ConstructSolution(Pheromone, ProblemSize, β, q0);
    Sicost ← Cost(Si);
    if Sicost ≤ Pbestcost then
      Pbestcost ← Sicost;
      Pbest ← Si;
    end
    LocalUpdateAndDecayPheromone(Pheromone, Si, Sicost, ρ);
  end
  GlobalUpdateAndDecayPheromone(Pheromone, Pbest, Pbestcost, ρ);
  while isUserInteraction() do
    GlobalAddAndRemovePheromone(Pheromone, Pbest, Pbestcost, ρ);
  end
end
return Pbest;

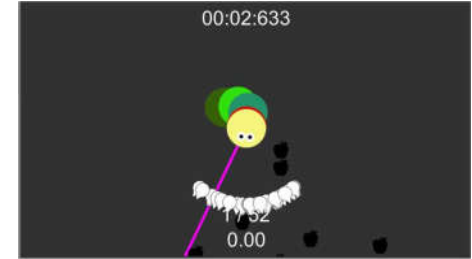
```

Holzinger, A., Plass, M., Holzinger, K., Crisan, G., Pintea, C. & Palade, V. 2016. Towards interactive Machine Learning (IML): Applying Ant Colony Algorithms to solve the Traveling Salesman Problem with the Human-in-the-Loop approach. Springer Lecture Notes in Computer Science LNCS 9817. 81-95, doi:10.1007/978-3-319-45507-56.

$$p_{ij} = \frac{[\tau_{ij}]^\alpha \cdot [\eta_{ij}]^\beta}{\sum_{l \in J_i^k} [\tau_{il}]^\alpha \cdot [\eta_{il}]^\beta}$$

- $p_{ij}$  ... **probability** of ants that they, at a particular node  $i$ , select the route from node  $i \rightarrow j$  (“**heuristic desirability**”)
- $\alpha > 0$  and  $\beta > 0$  ... the **influence parameters** ( $\alpha$  ... history coefficient,  $\beta$  ... heuristic coefficient) usually  $\alpha \approx \beta \approx 2 < 5$
- $\tau_{ij}$  ... the **pheromone value** for the components, i.e. the amount of pheromone on edge  $(i, j)$
- $k$  ... the set of usable components
- $J_i$  ... the set of nodes that ant  $k$  can reach from  $v_i$  (tabu list)
- $\eta_{ij} = \frac{1}{d_{ij}}$  ... attractiveness computed by a heuristic, indicating the “a-priori **desirability**” of the move

http://hci-kdd.org/gamification-interactive-machine-learning/



## LIVE DEMO

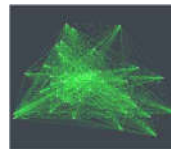
(<https://iml.hci-kdd.org/imlTspSolver/>)

ANDROID:

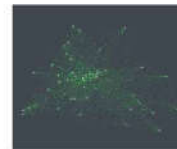
<https://play.google.com/store/apps/details?id=com.hcikdd.implacosolver>



- The pheromones are showing “the state” (high or low frequented paths of ants) of the algorithm.



initial pheromone distribution

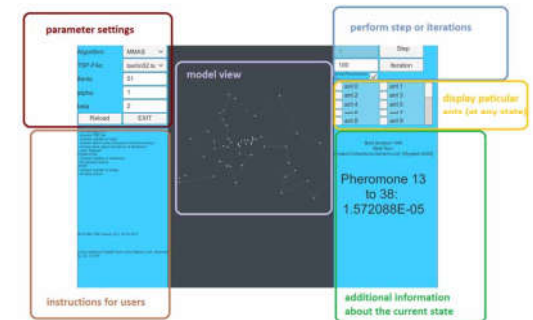


pheromones after 100 iterations

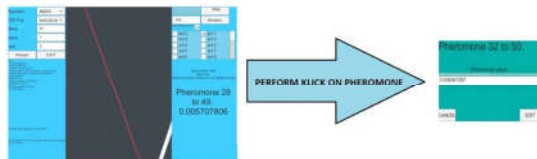


pheromones after 500 iterations

<http://iml.hci-kdd.org/imlTspSolver/>



- iteration vs. step: look inside the iteration
- make the ant algorithm interactive
  - change pheromones at any time
  - *change routes of certain ants in the current iteration (future work)*



## 06 Causality vs. Causability

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>

Causation is a matter of perception

We remember seeing the flame, and feeling a sensation called heat; without further ceremony, we call the one cause and the other effect

David Hume (1711-1776)

Statistical ML

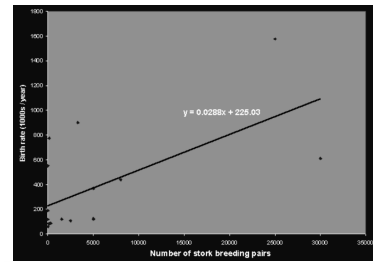
Forget causation! Correlation is all you should ask for.

Karl Pearson (1857-1936)

A mathematical definition of causality

Forget empirical observations! Define causality based on a network of known, physical, causal relationships

Judea Pearl (1936-)



Storks Deliver Babies ( $p = 0.008$ )

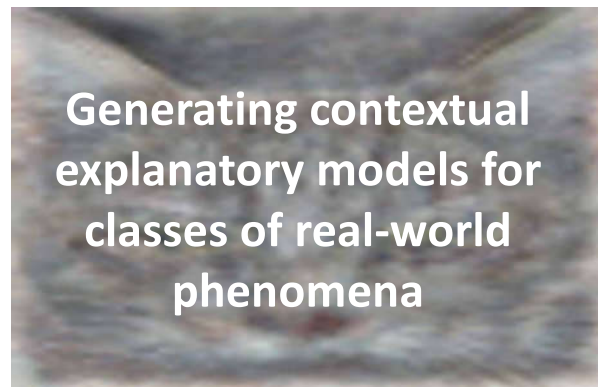
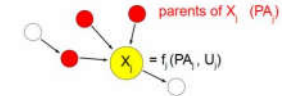
**KEYWORDS:**  
 Training:  
 Correlation:  
 Significance:  
 p-values:

**Robert Matthews**  
 Aston University, Birmingham, England.  
 e-mail: rajm@compuserve.com

**Summary**  
 This article shows that a highly statistically significant correlation exists between stork populations and human birth rates across Europe. While storks may not deliver babies, unthinking interpretation of correlation and p-values can certainly deliver unreliable conclusions.

### Functional Causal Model (Pearl et al.)

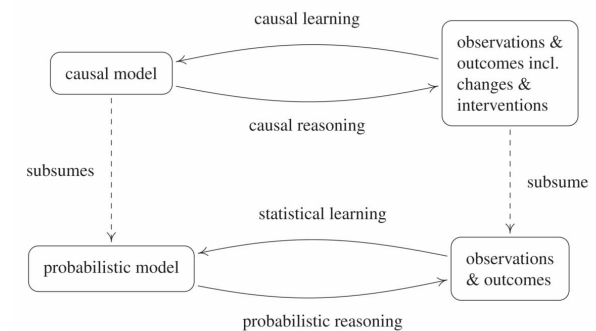
- Set of observables  $X_1, \dots, X_n$
- directed acyclic graph  $G$  with vertices  $X_1, \dots, X_n$
- Semantics: parents = direct causes
- $X_i = f_i(\text{ParentsOf}_i, \text{Noise}_i)$ , with independent  $\text{Noise}_1, \dots, \text{Noise}_n$ .
- "Noise" means "unexplained" (or "exogenous"), we use  $U_i$
- Can add requirement that  $f_1, \dots, f_n, \text{Noise}_1, \dots, \text{Noise}_n$  "independent" (cf. Lemeire & Dirks 2006, Janzing & Schölkopf 2010 — more below)



Explainability	in a technical sense highlights decision-relevant parts of the used representations of the algorithms and active parts in the algorithmic model, that either contribute to the model accuracy on the training set, or to a specific prediction for one particular observation. It does not refer to an explicit human model.
Causability	as the extent to which an explanation of a statement to a human expert achieves a specified level of causal understanding with effectiveness, efficiency and satisfaction in a specified context of use.

- Causability** := a property of a person, while
- Explainability** := a property of a system

## 07 explainable AI



Remember:  
Context !!!

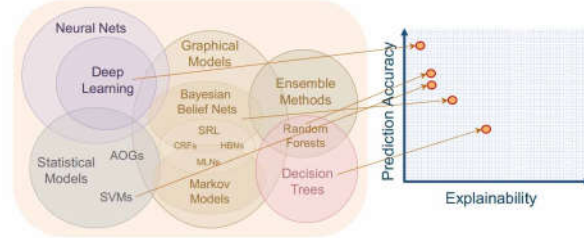




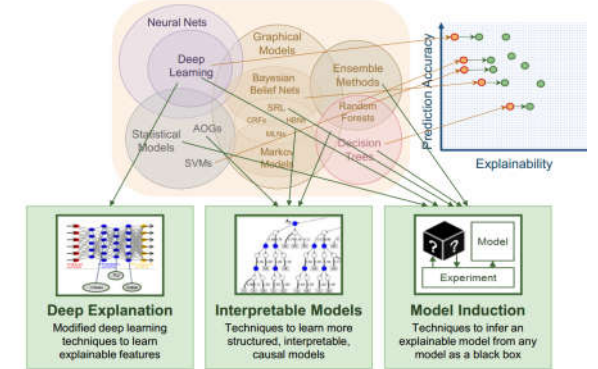
a woman riding a horse on a dirt road  
an airplane is parked on the tarmac at an airport  
a group of people standing on top of a beach

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 dee learning : [github.com/karpathy/neuraltalk2](https://github.com/karpathy/neuraltalk2)

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

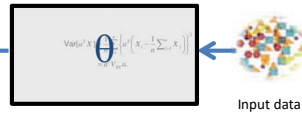


David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.

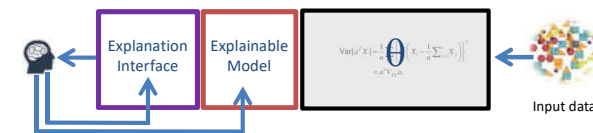


David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.

Why did the algorithm do that?  
Can I trust these results?  
How can I correct an error?



A possible solution

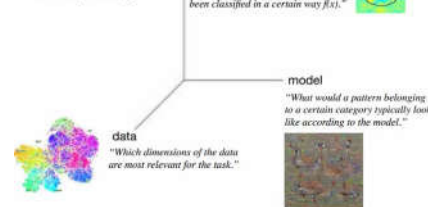


The domain expert can understand why ...  
The domain expert can learn and correct errors ...  
The domain expert can re-enact on demand ...

Post-hoc: Select a model and develop a technique to make it transparent



Different dimensions of "interpretability"



Ante-hoc: Select a model that is already transparent and optimize it

$$f(x) = \sum_{i=1}^d g_i(x_i)$$

- 1) Gradients
- 2) Sensitivity Analysis
- 3) Decomposition Relevance Propagation (Pixel-RP, Layer-RP, Deep Taylor Decomposition, ...)
- 4) Optimization (Local-IME – model agnostic, BETA transparent approximation, ...)
- 5) Deconvolution and Guided Backpropagation
- 6) Model Understanding
  - Feature visualization, Inverting CNN
  - Qualitative Testing with Concept Activation Vectors TCAV
  - Network Dissection

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

## Conclusion and Future Outlook

### Multi-Task Learning (MUTL)

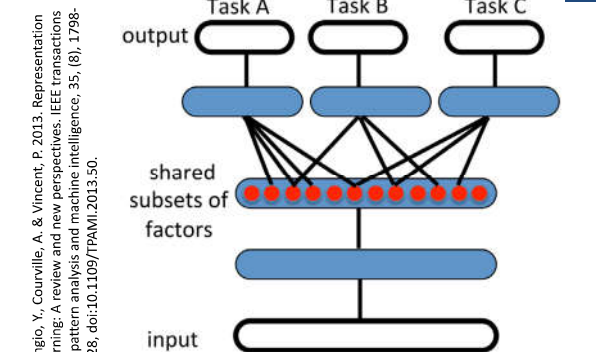
for improving prediction performance, help to reduce catastrophic forgetting

### Transfer learning (TRAL)

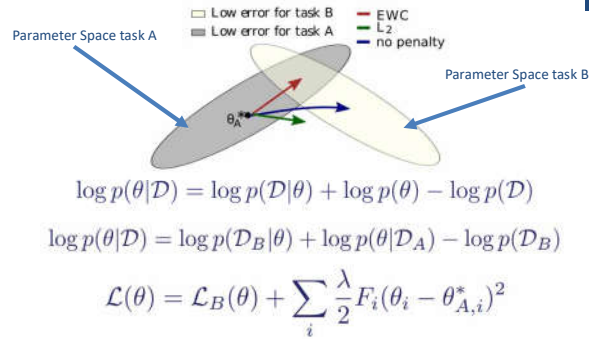
is not easy: learning to perform a task by exploiting knowledge acquired when solving previous tasks:  
a solution to this problem would have major impact to AI research generally and ML specifically.

### Multi-Agent-Hybrid Systems (MAHS)

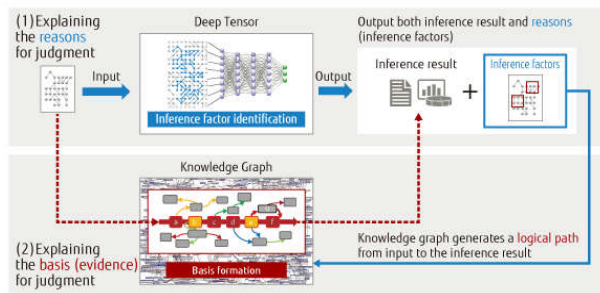
To include collective intelligence and crowdsourcing and making use of discrete models – avoiding to seek perfect solutions – better have a good solution < 5 min.



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.



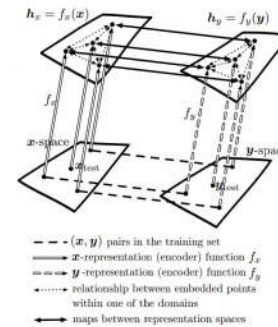
Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D. & Hadsell, R. 2016. Overcoming catastrophic forgetting in neural networks. arXiv preprint arXiv:1612.00796.



Randy Goebel, Ajay Chander, Katharina Holzinger, Freddy Lecue, Zeynep Akata, Simone Stumpf, Peter Kieseberg & Andreas Holzinger 2018. Explainable AI: the new 42? Springer Lecture Notes in Computer Science LNCS 11015

## Questions

- $x$  and  $y$  represent different modalities, e.g. text, sound, images, ...
- Generalization to new categories
- Larochelle et al. (2008) AAAI



Goodfellow, I., Bengio, Y. & Courville, A. 2016. Deep Learning, Cambridge: MIT Press, p.542

- Computers are fast, accurate and stupid,
- humans are 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>



- What is the HCI-KDD approach?
- What is meant by “integrative ML”?
- Why is a direct integration of AI-solutions into the workflow important?
- What are features?
- Why is understanding intelligence important?
- Why is understanding context even more important?
- What are currently the “best” ML-algorithms?
- What is the difference between Humanoid AI and Human-Level AI?
- Why is the health domain probably the most complex application domain for machine learning?

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



- Why are we speaking about “two different worlds” in the medical domain?
- Where is the problem in building the bridge between those two worlds?
- Why is the work of Bayes so important for machine learning?
- Why are Newton/Leibniz, Bayes/Laplace and Gauss so important for machine learning?
- What is learning and inference?
- What is the inverse probability?
- How does Bayesian optimization in principle work?



- What is the definition of aML?
- What is the best practice of aML?
- Why is “big data” necessary for aML?
- Provide examples for rare events!
- Give examples for NP-hard problems relevant for health informatics!
- Give the definition of iML?
- What is the benefit of a “human-in-the-loop”?
- Explain the differences of iML in contrast to supervised and semi-supervised learning!

- Active Learning
- Bayesian inference, Bayesian Learning
- Gaussian Processes
- Graphical Models
- Multi-Task Learning
- Reinforcement Learning
- Statistical Learning
- Transfer Learning
- Multi-Agent Hybrid Systems



## Multi-Task Feature Selection on Multiple Networks via Maximum Flows

Mahito Sugiyama<sup>1 (2)</sup>, Chloé-Agathe Azencott<sup>3</sup>, Dominik Grimm<sup>2,4</sup>, Yoshinobu Kawahara<sup>1</sup>, Karsten Borgwardt<sup>2,4</sup>

<sup>1</sup>Osaka University, <sup>2</sup>Max Planck Institutes Tübingen, <sup>3</sup>Mines ParisTech, Institut Curie, INSERM, <sup>4</sup>Eberhard Karls Universität Tübingen

Sugiyama, M., Azencott, C.-A., Grimm, D., Kawahara, Y. & Borgwardt, K. M. Multi-Task Feature Selection on Multiple Networks via Maximum Flows. SDM, 2014. 199-207.

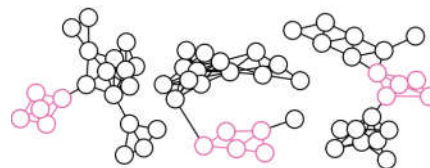
- What is causal relationship from purely observational data and why is it important?
- What is generalization?
- Why is understanding the context so important?
- What does the oracle in Active learning do?
- Explain catastrophic forgetting!
- Give an example for multi-task learning!
- What is the goal of transfer learning and why is this important for machine learning?
- Why would a contribution to a solution to transfer learning be a major breakthrough for artificial intelligence in general – and machine learning specifically?

- “The most interesting facts are those which can be used several times, those which have a chance of recurring ...*
- which, then, are the facts that have a chance of recurring?*
- In the first place, simple facts.”*



Henri Poincaré, Sciences et Methods (1908)

- Given multiple graphs
- Find features (=vertices), which are associated with the target response and tend to be connected to each other



# Appendix

- Bernhard Schölkopf (MPI Tübingen) <https://is.tuebingen.mpg.de/person/bs>
- Leslie Valiant (Harvard) <https://people.seas.harvard.edu/~valiant>
- Joshua Tenenbaum (MIT) <http://web.mit.edu/cocosci/josh.html>
- Andrew G. Wilson Cornell (Eric P. Xing, CMU) <https://people.orie.cornell.edu/andrew>
- Nando de Freitas (Oxford) <https://www.cs.ox.ac.uk/people/nando.defreitas>
- Yoshua Bengio (Montreal) [http://www.iro.umontreal.ca/~bengioy/yoshua\\_en](http://www.iro.umontreal.ca/~bengioy/yoshua_en)
- David Blei (Columbia) <http://www.cs.columbia.edu/~blei>
- Zoubin Ghahramani (Cambridge) <http://mlg.eng.cam.ac.uk/zoubin>
- Noah Goodman (Stanford) <http://cocolab.stanford.edu/ndg.html>

$$\underset{\substack{S_1, \dots, S_K \subseteq V \\ K \text{ tasks}}}{\operatorname{argmax}} \sum_{i=1}^K \left( \underbrace{f_i(S_i)}_{\text{association}} - g_i(S_i) \right) - \underbrace{\sum_{i < j} h(S_i, S_j)}_{\text{penalty}}$$

$$f_i(S_i) := \sum_{v \in S_i} q_i(v), \quad g_i(S_i) := \lambda \sum_{\substack{e \in E_i \\ \text{connectivity}}} w_i(e) + \underbrace{\eta |S_i|}_{\text{sparsity}}$$

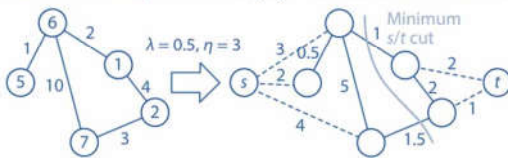
$$h(S_i, S_j) := \mu |S_i \Delta S_j| = \mu |(S \cup S') \setminus (S \cap S')|$$

- efficiently solved by max-flow algorithms
- performance is superior to Lasso-based methods

Sugiyama, M., Azencott, C.-A., Grimm, D., Kawahara, Y. & Borgwardt, K. M. Multi-Task Feature Selection on Multiple Networks via Maximum Flows. SDM, 2014. 199-207.

- Networks (graphs) are everywhere in health informatics
- Biological pathways (KEGG), chemical compounds, (PubChem), social networks, ...
- Question often: Which part of the network is responsible for performing a particular function?
- Feature selection on networks
  - Features = vertices (nodes)
  - Network topology = a priori knowledge of relationships between features
- Multi-task feature selection should be considered for more effectiveness**

- The  $s/t$ -network  $M(G) = (V \cup \{s, t\}, E \cup S \cup T)$  with  $S = \{\{s, v\} \mid v \in V, q(v) > \eta\}$ ,  $T = \{\{t, v\} \mid v \in V, q(v) < \eta\}$  and set the capacity  $c : E' \rightarrow \mathbb{R}^+$  to  $c(\{v, u\}) = \begin{cases} |q(u) - \eta| & \text{if } u \in \{s, t\} \text{ and } v \in V, \\ \lambda w(\{v, u\}) & \text{otherwise} \end{cases}$
- The minimum  $s/t$  cut of  $M(G)$  = the solution of SConES



Azencott, C.-A., Grimm, D., Sugiyama, M., Kawahara, Y. & Borgwardt, K. M. 2013. Efficient network-guided multi-locus association mapping with graph cuts. *Bioinformatics*, 29, (13), i171-i179.

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

Frequency	PMI	related body balance	SCP		
joint	5967	swollen joints	5.8	swollen joints	0.002
Arthritis	661	ASPIRIN INTOLERANCE	5.8	pain	0.001
fatigue	429	Epinephrine hypokalemia	5.8	Arthritis	0.001
Arthritis	301	swollen joints	5.4	fatigue	0.0007
swollen joints	299	Joint tenderness	5.1	swollen	0.0007
erythema	255	Occipital headache	4.2	ophthalmology	0.0004
Back Pain	254	Neuromuscular excitation	4.2	Back Pain	0.0004
headache	239	Reaction sleep	5.8	polydipsia	0.0004
ophthalmology	228	joint capsule	5.7	joint stiffness	0.0004
Anorexia	221	joint symptoms	5.3	Joint tenderness	0.0004
dysequia	218	Partial face	5.3	leg pain	0.0004
weakness	210	feeling of malaise	4.5	metastasis	0.0004
nausea	199	Morvan's sign	5.4	skin manifestations	0.0007
Reactivity of Function	181	diffuse pain	4.3	neck pain	0.0007
low back pain	167	Palmar erythema	5.3	Eye Manifestations	0.0004
abdominal pain	141	Abnormal sensation	5.2	low back pain	0.0004

Holzinger, A., Yildirim, P., Geier, M. & Simon, K.-M. 2013. Quality-Based Knowledge Discovery from Medical Text on the Web. In: Pasi, G., Bordogna, G. & Jain, L. C. (eds.) *Quality Issues in the Management of Web Information*, Intelligent Systems Reference Library, ISRL 50. Berlin Heidelberg: Springer, pp. 145-158, doi:10.1007/978-3-642-37688-7\_7.

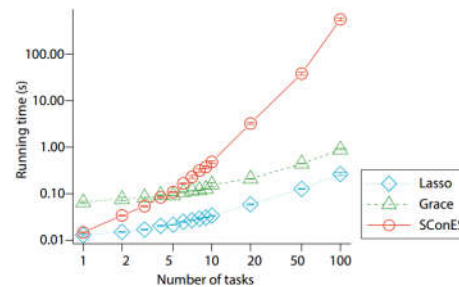
- Single task feature selection on a network
- Given a weighted graph  $G = (V, E)$ 
  - Each  $v \in V$  has a relevance score  $q(v)$
  - If you have a design matrix  $\mathbf{X} \in \mathbb{R}^{N \times |V|}$
  - and a response vector  $\mathbf{y} \in \mathbb{R}^N$   $\#(v)$  is the association of  $\mathbf{y}$  and each feature of  $\mathbf{X}$

Goal: Find a subset  $S \subset V$  which maximizes

$$f(S) := \sum_{v \in S} q(v)$$

while  $S$  is small and vertices are connected

Azencott, C.-A., Grimm, D., Sugiyama, M., Kawahara, Y. & Borgwardt, K. M. 2013. Efficient network-guided multi-locus association mapping with graph cuts. *Bioinformatics*, 29, (13), i171-i179.



Azencott, C.-A., Grimm, D., Sugiyama, M., Kawahara, Y. & Borgwardt, K. M. 2013. Efficient network-guided multi-locus association mapping with graph cuts. *Bioinformatics*, 29, (13), i171-i179.

$$\begin{aligned} & \cdot \operatorname{argmax}_{S \subset V} f(S) - g(S) \\ & f(S) := \sum_{v \in S} q(v), \quad g(S) := \underbrace{\lambda \sum_{e \in B} w(e)}_{\text{connectivity}} + \underbrace{\eta |S|}_{\text{sparsity}} \\ & - B = \{\{v, u\} \in E \mid v \in V \setminus S, u \in S\} \text{ (boundary)} \\ & - w : E \rightarrow \mathbb{R}^+ \text{ is a weighting function} \end{aligned}$$



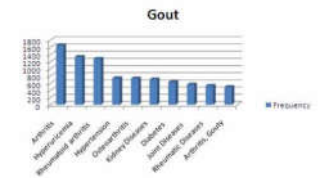
Azencott, C.-A., Grimm, D., Sugiyama, M., Kawahara, Y. & Borgwardt, K. M. 2013. Efficient network-guided multi-locus association mapping with graph cuts. *Bioinformatics*, 29, (13), i171-i179.

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

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

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

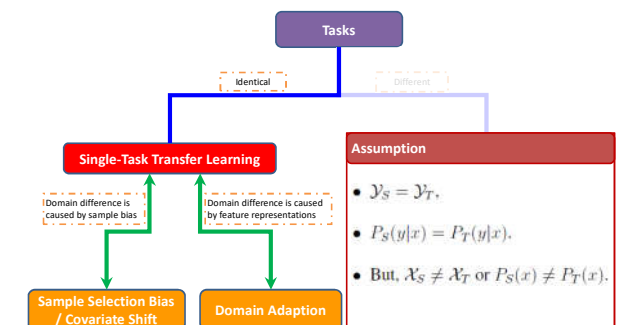
$$\begin{aligned} P(w_i) &= \frac{94,834}{20,033,079} & P(w_j) &= \frac{74}{20,033,079} \\ P(w_i, w_j) &= \frac{13}{94,834} & PMI(w_i, w_j) &= 7.7. \end{aligned}$$



Holzinger, A., Simon, K. M. & Yildirim, P. Disease-Disease Relationships for Rheumatic Diseases: Web-Based Biomedical Textmining an Knowledge Discovery to Assist Medical Decision Making. 36th Annual IEEE Computer Software and Applications Conference (COMPSAC), 16-20 July 2012 1212 Izmir: IEEE, 573-580, doi:10.1109/COMPSAC.2012.77.

- Motivation: If two domains are related to each other, then there may exist some “pivot” features across both domain.
- Pivot features are features that behave in the same way for discriminative learning in both domains.
- Main Idea: To identify correspondences among features from different domains by modeling their correlations with pivot features.
- Non-pivot features form different domains that are correlated with many of the same pivot features are assumed to correspond, and they are treated similarly in a discriminative learner.
- Blitzer, J., McDonald, R. & Pereira, F. Domain adaptation with structural correspondence learning. *Proceedings of the 2006 conference on empirical methods in natural language processing*, 2006. Association for Computational Linguistics, 120-128.

Blitzer, J., McDonald, R. & Pereira, F. Domain adaptation with structural correspondence learning. *Proceedings of the 2006 conference on empirical methods in natural language processing*, 2006. Association for Computational Linguistics, 120-128.





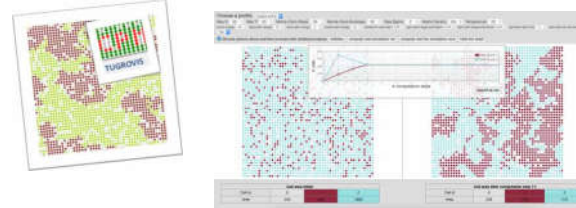
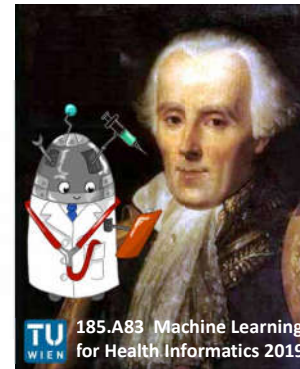
# Open Problem: How to avoid negative transfer?

- Example 1: Subspace Clustering
- Example 2: k-Anonymization
- Example 3: Protein Design

Hund, M., Böhm, D., Sturm, W., Sedlmair, M., Schreck, T., Ullrich, T., Keim, D. A., Majnarić, L. & Holzinger, A. 2016. Visual analytics for concept exploration in subspaces of patient groups: Making sense of complex datasets with the Doctor-in-the-loop. Brain Informatics, 1-15, doi:10.1007/s40708-016-0043-5.

Kieseberg, P., Malle, B., Fruehwirt, P., Weippl, E. & Holzinger, A. 2016. A tamper-proof audit and control system for the doctor in the loop. Brain Informatics, 3, (4), 269-279, doi:10.1007/s40708-016-0046-2.

Lee, S. & Holzinger, A. 2016. Knowledge Discovery from Complex High Dimensional Data. In: Michaelis, S., Piatkowski, N. & Stolpe, M. (eds.) Solving Large Scale Learning Tasks. Challenges and Algorithms, Lecture Notes in Artificial Intelligence LNAI 9580. Springer, pp. 148-167, doi:10.1007/978-3-319-41706-6\_7.

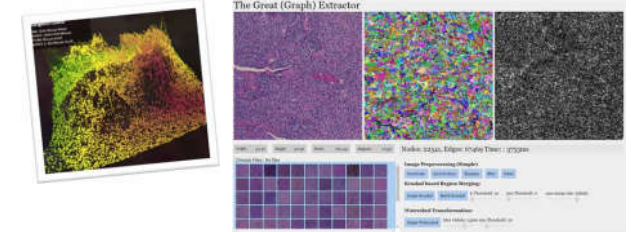


- Contribute to understanding tumor growth
- Goal: Help to Refine → Reduce → Replace
- Towards discrete Multi-Agent Hybrid Systems

Jeanquartier, F., Jean-Quartier, C., Cemernek, D. & Holzinger, A. 2016. In silico modeling for tumor growth visualization. BMC Systems Biology, 10, (1), 1-15, doi:10.1186/s12918-016-0318-8.

Jeanquartier, F., Jean-Quartier, C., Kotlyar, M., Tokar, T., Hauschild, A.-C., Jurisica, I. & Holzinger, A. 2016. Machine Learning for In Silico Modeling of Tumor Growth. In: Springer Lecture Notes in Artificial Intelligence LNAI 9605. Cham: Springer International Publishing, pp. 415-434, doi:10.1007/978-3-319-50478-0\_21.

- Computational resource intensive (supercomps, cloud CPUs, **federated learning**, ...)
- Black-Box approaches – lack **transparency**, do not foster trust and acceptance among end-user, legal aspects make “black box” difficult!
- **Non-convex**: difficult to set up, to train, to optimize, needs a lot of expertise, error prone
- Very bad in dealing with **uncertainty**
- **Data intensive**, needs often millions of training samples ...



- Contribute to graph understanding and algorithm prototyping by real-time visualization, interaction and manipulation
- Supports client-based federated learning
- Towards an online graph exploration and analysis platform

Malle, B., Kieseberg, P., Weippl, E. & Holzinger, A. 2016. The right to be forgotten: Towards Machine Learning on perturbed knowledge bases. Springer Lecture Notes in Computer Science LNCS 9817. Heidelberg, Berlin, New York: Springer, pp. 251-256, doi:10.1007/978-3-319-45507-5\_17.