



Trustworthy Learning and Reasoning in Complex Domains

Federico Cerutti — federico.cerutti@unibs.it

Augmenting human sensemaking abilities
to achieve causal insights and foresights

(a.k.a. *situational understanding*)



Overture. A brief historical case.

Act I. On conjectures, refutations, and argumentation.

Act II. There is no certain datum in the world.

Act III. Interesting problems are complex.

Epilogue.

the use of the intellectual faculty; to comprehend; to be informed by another; to learn. **understanding**, un-dér-stand'ing, *n.* Intelligent; knowing; skilful. — *n.* The act of one who understands; comprehension; apprehension; discernment; knowledge; **clear insight**; the faculty or power by which one understands; the faculty of the human mind otherwise known as the intellect; **the power of thinking and reasoning**; intelligence between two or more persons; agreement of minds; anything mutually understood or agreed upon. **understate**, un-dér-stát', *v.t.* To state too low; to state or represent less strongly than the truth will bear. **understatement**, un-dér-stát'ment, *n.* understating; a statement under

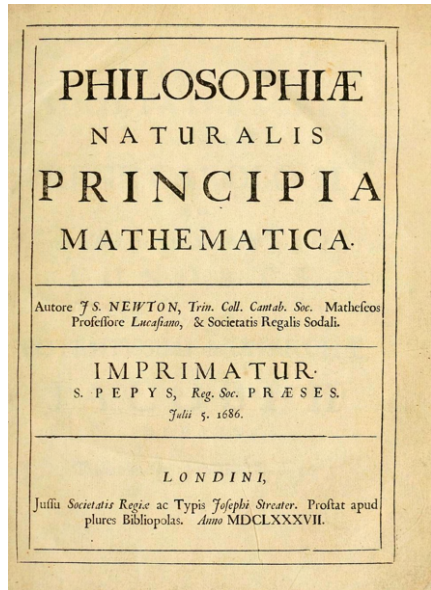
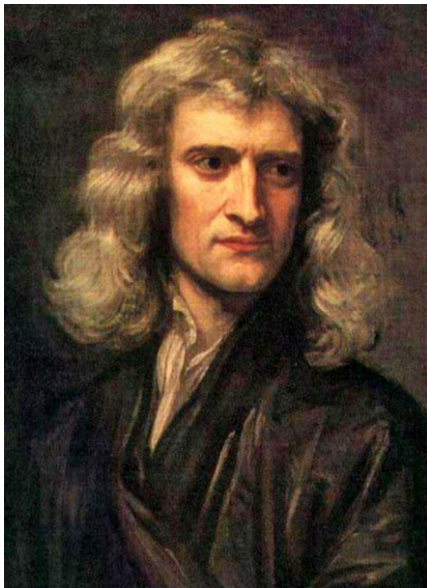
determine
not certain
undeterm
not restr
undevia
viating;
ciple, or
undiges
by the st
arranged
undigni
fied; sho
undilut
or mixed
any adm
underline
A water

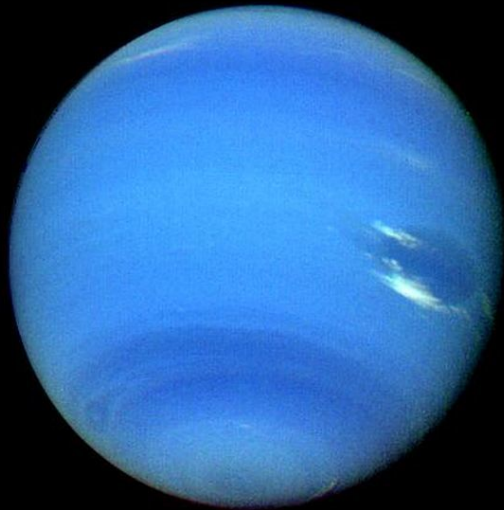
Empiricism

All hypotheses and theories must be tested against observations of the natural world, rather than resting solely on a priori reasoning, intuition, or revelation.









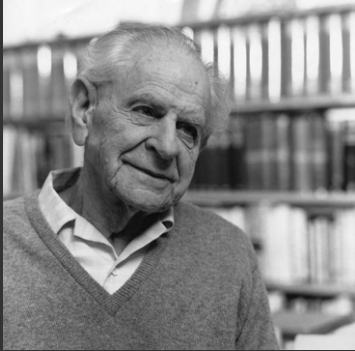




The path of the planet Uranus did not conform to the path predicted by Newton's law of gravitation in presence of the known planets.

Explanations:

- Human/instrument measure error
- Newton's laws are mistaken
- An invisible magic teapot caused the perturbation in order to show the *hubris* of modern science
- ...
- Newton's laws—confirmed by a significant amount of evidence—are correct and the perturbation is caused by another, unknown, planet



Scientific theories are capable of being refuted: they are **falsifiable**

Verification and falsification are different processes:

- No accumulation of confirming instances is sufficient
- Only one contradicting instance suffices to refute a theory

Scientific theories are tentative

Overture. A brief historical case.

Act I. On conjectures, refutations, and argumentation.

Act II. There is no certain datum in the world.

Act III. Interesting problems are complex.

Epilogue.

Does MMR vaccination cause autism?

Argument from Correlation to Cause

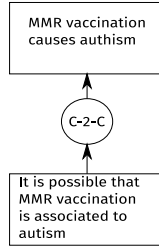
Correlation Premise: There is a positive correlation between A and B.

Conclusion: A causes B.

CQ1: Is there really a correlation between A and B?

CQ2: Is there any reason to think that the correlation is any more than a coincidence?

CQ3: Could there be some third factor, C, that is causing both A and B?



Early report

Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children

A J Wakefield, S H Murch, A Anthony, J Linnell, D M Casson, M Malik, M Berelowitz, A P Dhillon, M A Thomson, P Harvey, A Valentine, S E Davies, J A Walker-Smith

Summary

We investigated a consecutive

Introduction

We saw several children who, after a period of

Findings Onset of behavioural symptoms was associated, by the parents, with measles, mumps, and rubella vaccination in eight of the 12 children, with measles infection in one child, and otitis media in another. All 12 children had intestinal abnormalities, ranging from lymphoid nodular hyperplasia to aphthoid ulceration. Histology showed patchy chronic inflammation in the colon in 11 children and reactive ileal lymphoid hyperplasia in seven, but no granulomas. Behavioural disorders included autism (nine), disintegrative psychosis (one), and possible postviral or vaccinal encephalitis (two). There were no focal neurological abnormalities and MRI and EEG tests were normal. Abnormal laboratory results were significantly raised urinary methylmalonic acid compared with age-matched controls ($p=0.003$), low haemoglobin in four children, and a low serum IgA in four children.

Support

Child	Behavioural diagnosis	Exposure identified by parents or doctor	Interval from exposure to first behavioural symptom	Features associated with exposure	Age at onset of first symptom	
					Behaviour	Bowel
1	Autism	MMR	1 week	Fever/delirium	12 months	Not known
2	Autism	MMR	2 weeks	Self injury	13 months	20 months
3	Autism	MMR	48 h	Rash and fever	14 months	Not known
4	Autism?	MMR	Measles vaccine at 15 months followed by slowing in development. Dramatic deterioration in behaviour immediately after MMR at 4.5 years	Repetitive behaviour, self injury, loss of self-help	4.5 years	18 months
5	Disintegrative disorder?	None—MMR at 16 months	Self-injurious behaviour started at 18 months		4 years	
6	Autism	MMR	1 week		15 months	18 months
7	Autism	MMR	24 h	Rash & convulsion; gaze avoidance & self injury	15 months	18 months
8	Post-vaccinal encephalitis?	MMR	2 weeks	Convulsion, gaze avoidance	21 months	2 years
9	Autistic spectrum disorder	Recurrent otitis media	1 week (MMR 2 months previously)	Fever, convulsion, rash & diarrhoea	19 months	19 months
10	Post-viral encephalitis?	Measles (previously vaccinated with MMR)	24 h	Disinterest; lack of play	18 months	2.5 years
11	Autism	MMR	1 week	Fever, rash & vomiting	15 months	Not known
12	Autism	None—MMR at 15 months	Loss of speech development and deterioration in language skills noted at 16 months	Recurrent "viral pneumonia" for 8 weeks following MMR	15 months	Not known

MMR=measles, mumps, and rubella vaccine.

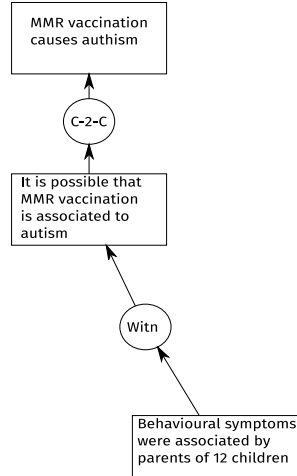
Table 2: Neuropsychiatric diagnosis

What else should be true if the causal link is true?

with autoimmune hepatitis.¹⁴

We did not prove an association between measles, mumps, and rubella vaccine and the syndrome described. Virological studies are underway that may help to resolve this issue.

If there is a causal link between measles, mumps, and rubella vaccine and this syndrome, a rising incidence might be anticipated after the introduction of this vaccine in the UK in 1988. Published evidence is inadequate to show whether there is a change in incidence¹⁵ or a link with measles, mumps, and rubella vaccine.¹⁶ A genetic predisposition to autistic-spectrum disorders is suggested by over-representation in boys and a greater concordance rate in monozygotic than in dizygotic twins.¹⁷ In the context of susceptibility to infection, a genetic association



The New England Journal of Medicine

Copyright © 2002 by the Massachusetts Medical Society

VOLUME 347

NOVEMBER 7, 2002

NUMBER 19



A POPULATION-BASED STUDY OF MEASLES, MUMPS, AND RUBELLA VACCINATION AND AUTISM

KREESTEN MELDGAARD MADSEN, M.D., ANDERS HVIID, M.Sc., MOGENS VESTERGAARD, M.D., DIANA SCHENDEL, Ph.D.,
JAN WOHLFAHRT, M.Sc., POUL THORSEN, M.D., JØRN OLSEN, M.D., AND MADS MELBYE, M.D.

ABSTRACT

that vaccina-

suggested that the measles

vaccine can

There was no association between the age at the time of vaccination, the time since vaccination, or the date of vaccination and the development of autistic disorder.

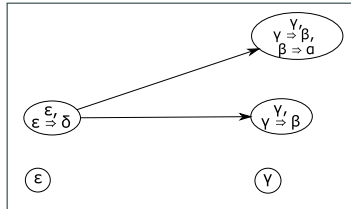
Conclusions This study provides strong evidence against the hypothesis that MMR vaccination causes autism. (N Engl J Med 2002;347:1477-82.)

Copyright © 2002 Massachusetts Medical Society.

Support



Results Of the 537,303 children in the cohort (representing 2,129,864 person-years), 440,655 (82.0 percent) had received the MMR vaccine. We identified 316 children with a diagnosis of autistic disorder and 422 with a diagnosis of other autistic-spectrum disorders. After adjustment for potential confounders, the relative risk of autistic disorder in the group of vaccinated children, as compared with the unvaccinated group, was 0.92 (95 percent confidence interval, 0.68 to 1.24), and the relative risk of another autistic-spectrum disorder was 0.83 (95 percent confidence interval, 0.65 to 1.07).

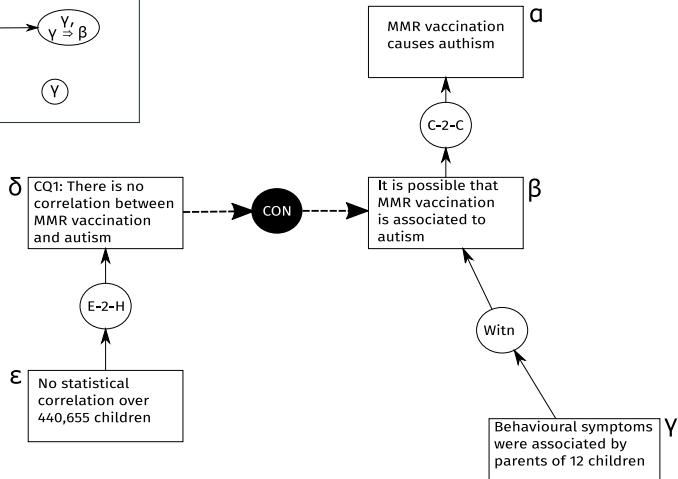


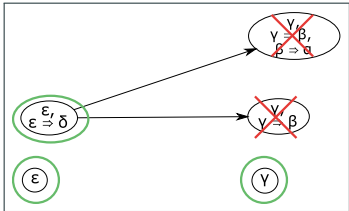
$$\beta \implies \alpha$$

$$\gamma \implies \beta$$

$$\epsilon \implies \delta$$

$$\delta \in \bar{\beta}$$



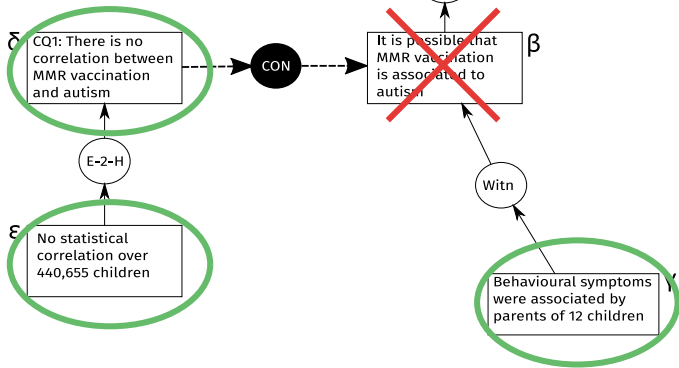


$$\beta \Rightarrow \alpha$$

$$\gamma \Rightarrow \beta$$

$$\epsilon \Rightarrow \delta$$

$$\delta \in \bar{\beta}$$



Results (tiny summary)

HCI Assessment of argumentation semantics against human intuition (ECAI 2014)

Algorithms Efficient algorithms and ensemble approaches (KR 2014, AAI 2015, ECAI 2016, KER 2018, IJAR 2018, AIJ 2019, IJCAI 2021)

Impact Implementation in the CIsaces.org online system (AAMAS 2015, SPIE 2018, COMMA 2018, JURIX 2018, AI³ 2021)

Fact extraction from Twitter

Extract

rt @breakingnews rumors of nyse trading floor rioting are not true says nyse

Text

RT @BreakingNews: Rumors of NYSE trading floor rioting are not true, says NYSE - @politico @CNBC @weatherchannel

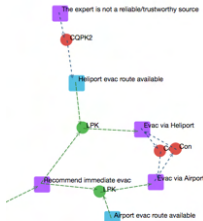
Twitter URI

<https://twitter.com/LasiewickiAnn/status/263222115120082945>

Time

Thu Nov 01 2012 10:13:37 GMT+0000 (GMT)

Argumentation graph manipulation



Natural Language Generation for Automatic Reporting

Report

We have reasons to believe that:

- This claim is not supported by evidence
- Recommend immediate evac because RISK TO LIFE
- RISK TO LIFE because UK nationals at NYU hospital, and [info received] nyu hospital still being evacuated rioting and fires
- UK nationals at NYU hospital because [info received] Embassy report UK nationals at NYU hospital

Moreover, we also have the following 2 hypotheses.

Hypothesis number 1

- Evac via Airport because Recommend immediate evac, and [info received] Airport evac route available

Hypothesis number 2

- Evac via Heliport because Recommend immediate evac, and [info received] Heliport evac route available

Here the pieces of information we received

- Airport evac route available
- Embassy report UK nationals at NYU hospital
- Reports of riots confirmed
- Heliport evac route available
- nyu hospital still being evacuated rioting and fires

Available for use by professional analysts in the US Army Research Laboratory, and the UK Joint Forces Intelligence Group

TRL4: validation in a laboratory environment

<https://tiresia.unibs.it/cispaces>

Overture. A brief historical case.

Act I. On conjectures, refutations, and argumentation.

Act II. There is no certain datum in the world.

Act III. Interesting problems are complex.

Epilogue.

Qualification problem

“ For example, the successful use of a boat to cross a river requires, if the boat is a rowboat, that the oars and rowlocks be present and unbroken, and that they fit each other. Many other qualifications can be added, making the rules for using a rowboat almost impossible to apply, and yet anyone will still be able to think of additional requirements not yet stated.

”

J. McCarthy, “Circumscription—A Form of Nonmonotonic Reasoning,” *AIJ*, 13 (12): 2739, 1980.

Uncertainty

Reliability of the Source

- A** Completely reliable
- B** Usually reliable
- C** Fairly reliable
- D** Not usually reliable
- E** Unreliable
- F** Reliability cannot be judged

Credibility of the Information

- 1** Confirmed by other sources
- 2** Probably true
- 3** Possibly true
- 4** Doubtful
- 5** Improbable
- 6** Truth cannot be judged


```

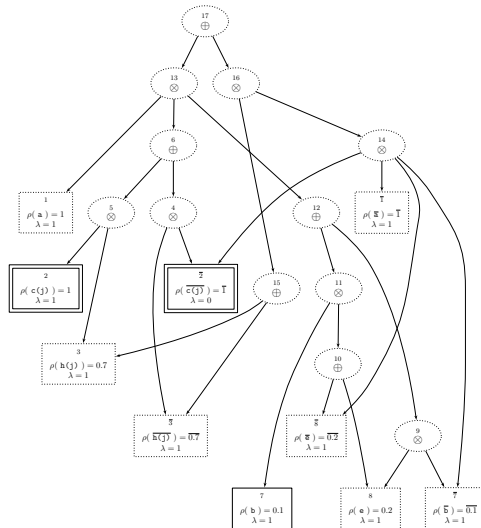
0.1:: burglary.
0.2:: earthquake.
0.7:: hears_alarm(john).
alarm :- burglary.
alarm :- earthquake.
calls(john) :- alarm, hears_alarm(john).
evidence(calls(john)).
query(burglary).

```

```

alarm  $\leftrightarrow$  burglary  $\vee$  earthquake
calls(john)  $\leftrightarrow$  alarm  $\wedge$  hears_alarm(john)
calls(john)

```



Where numbers come from?

# Day	Earthquake
1	T
2	T
3	F
4	F
5	F
6	F
7	F
8	F
9	F
10	F

π : true—unknown—probability of earthquake in a given period of time

Let y be the number of occurrence of earthquake per period of time ($y = 2$)

From Bayes' theorem, we can estimate the **posterior** distribution of π given the data on the basis of a **prior**: $g(\pi|y) \propto g(\pi) \cdot f(y|\pi)$

The conjugate of a binomial is the Beta distribution. If:

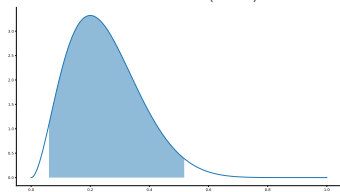
$$g(\pi; a, b) = \text{Beta}(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \pi^{a-1} (1-\pi)^{b-1}$$

then: $g(\pi|y) = \text{Beta}(y+a, n-y+b)$

If $a = b = 1$ (uniform prior), then $g(\pi|y) = \text{Beta}(y+1, n-y+1)$

In the example, $g(\pi|y=2, n=10) = \text{Beta}(3, 9)$

$$X_1 \sim \text{Beta}(3, 9)$$

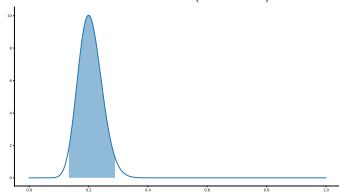


$$E[X_1] = 0.2500$$

$$\text{Var}(X_1) = 1.4423 \cdot 10^{-2}$$

95% Confidence Interval:
[0.0602, 0.5178]

$$X_2 \sim \text{Beta}(21, 81)$$

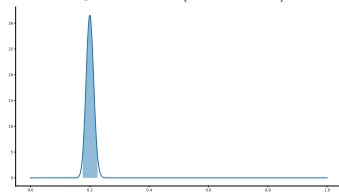


$$E[X_2] = 0.2059$$

$$\text{Var}(X_2) = 1.5873 \cdot 10^{-3}$$

95% Confidence Interval:
[0.1336, 0.2891]

$$X_3 \sim \text{Beta}(201, 801)$$



$$E[X_3] = 0.2006$$

$$\text{Var}(X_3) = 1.5988 \cdot 10^{-4}$$

95% Confidence Interval:
[0.1764, 0.2259]

Although $E[X_1] \simeq E[X_2] \simeq E[X_3] \simeq 0.2$
they represent remarkably different random variables

Microsoft Human-AI Interaction Guidelines

Guideline 1: Make clear what the system can do.

Guideline 2: Make clear how well the system can do what it can do.

...

S. Amershi et. al., "Guidelines for Human-AI Interaction," CHI 2019

EU Requirements of Trustworthy AI

Human agency and oversight

Technical robustness and safety

Privacy and data governance

Transparency

Diversity, non-discrimination, and fairness

Societal and environmental wellbeing

Accountability

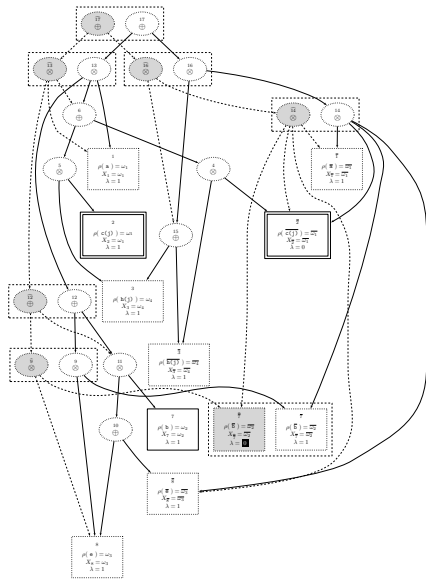
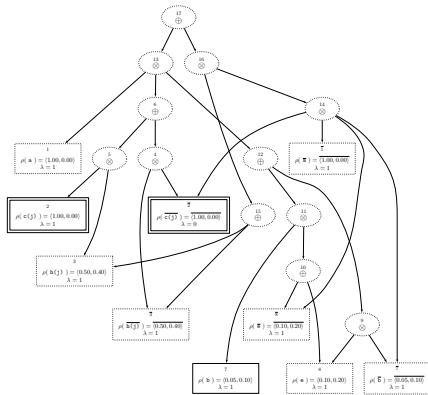
EUROPEAN COMMISSION, 2019. High-Level Expert Group on Artificial Intelligence.

```

 $\omega_2$  :: burglary .
 $\omega_3$  :: earthquake .
 $\omega_4$  :: hears_alarm (john) .
alarm :- burglary .
alarm :- earthquake .
calls (john) :- alarm , hears_alarm (john) .
evidence (calls (john)) .
query (burglary) .

```

Identifier	Beta parameters
ω_1	Beta(∞ , 1)
$\overline{\omega_1}$	Beta(1, ∞)
ω_2	Beta(2, 18)
$\overline{\omega_2}$	Beta(18, 2)
ω_3	Beta(2, 8)
$\overline{\omega_3}$	Beta(8, 2)
ω_4	Beta(3.5, 1.5)
$\overline{\omega_4}$	Beta(1.5, 3.5)



Cerutti, Kaplan, Kimmig, Şensoy, Handling Epistemic and Aleatory Uncertainties in Probabilistic Circuits, Under Submission, 2021, <https://arxiv.org/abs/2102.10865>

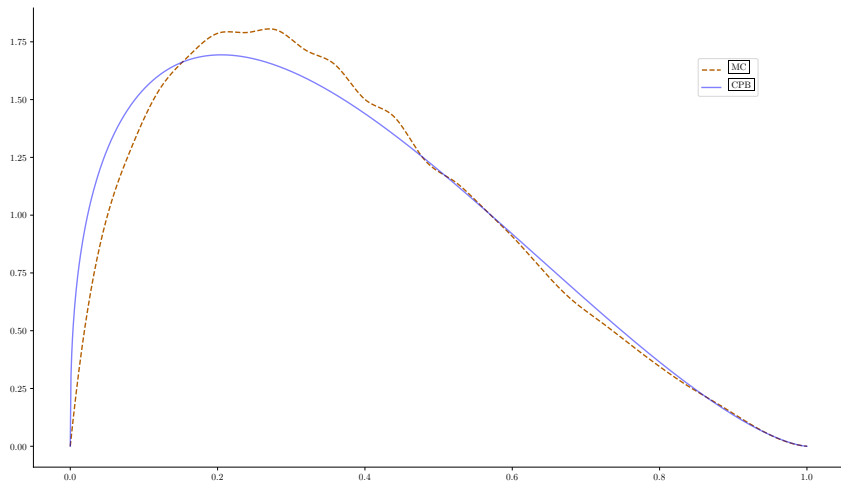
Let n be a \oplus -gate over C nodes, its children

$$\begin{aligned}\mathbb{E}[X_n] &= \sum_{c \in C} \mathbb{E}[X_c], \\ \text{cov}[X_n] &= \sum_{c \in C} \sum_{c' \in C} \text{cov}[X_c, X_{c'}], \\ \text{cov}[X_n, X_z] &= \sum_{c \in C} \text{cov}[X_c, X_z] \quad \text{for } z \in \widehat{N}_A \setminus \{n\}\end{aligned}$$

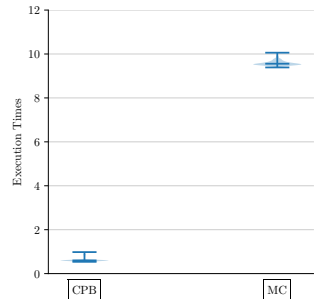
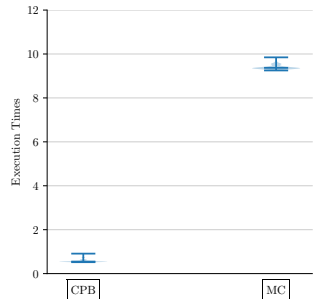
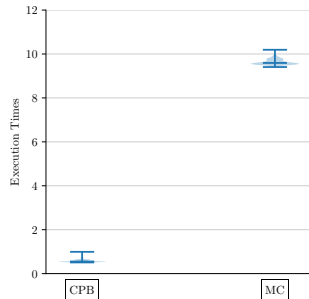
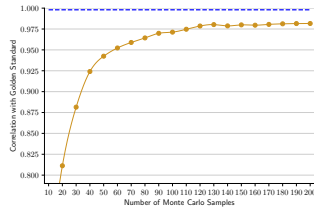
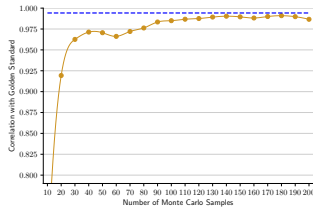
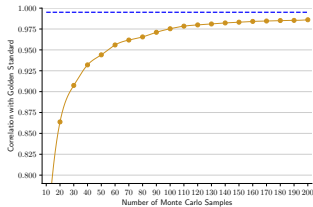
Let n be a \otimes -gate over C nodes, its children

$$\begin{aligned}\mathbb{E}[X_n] &= \prod_{c \in C} \mathbb{E}[X_c], \\ \text{cov}[X_n] &\simeq \sum_{c \in C} \sum_{c' \in C} \frac{\mathbb{E}[X_n]^2}{\mathbb{E}[X_c]\mathbb{E}[X_{c'}]} \text{cov}[X_c, X_{c'}], \\ \text{cov}[X_n, X_z] &\simeq \sum_{c \in C} \frac{\mathbb{E}[X_n]}{\mathbb{E}[X_c]} \text{cov}[X_c, X_z] \quad \text{for } z \in \widehat{N}_A \setminus \{n\}.\end{aligned}$$

$$\begin{aligned}\mathbb{E}\left[\frac{X_r}{X_{\hat{r}}}\right] &\simeq \frac{\mathbb{E}[X_r]}{\mathbb{E}[X_{\hat{r}}]}, \\ \text{cov}\left[\frac{X_r}{X_{\hat{r}}}\right] &\simeq \frac{1}{\mathbb{E}[X_{\hat{r}}]^2} \text{cov}[X_r] + \frac{\mathbb{E}[X_r]^2}{\mathbb{E}[X_{\hat{r}}]^4} \text{cov}[X_{\hat{r}}] - 2 \frac{\mathbb{E}[X_r]}{\mathbb{E}[X_{\hat{r}}]^3} \text{cov}[X_r, X_{\hat{r}}].\end{aligned}$$



Cerutti, Kaplan, Kimmig, Şensoy, Handling Epistemic and Aleatory Uncertainties in Probabilistic Circuits, Under Submission, 2021, <https://arxiv.org/abs/2102.10865>



Cerutti, Kaplan, Kimmig, Şensoy, Handling Epistemic and Aleatory Uncertainties in Probabilistic Circuits, Under Submission, 2021, <https://arxiv.org/abs/2102.10865>

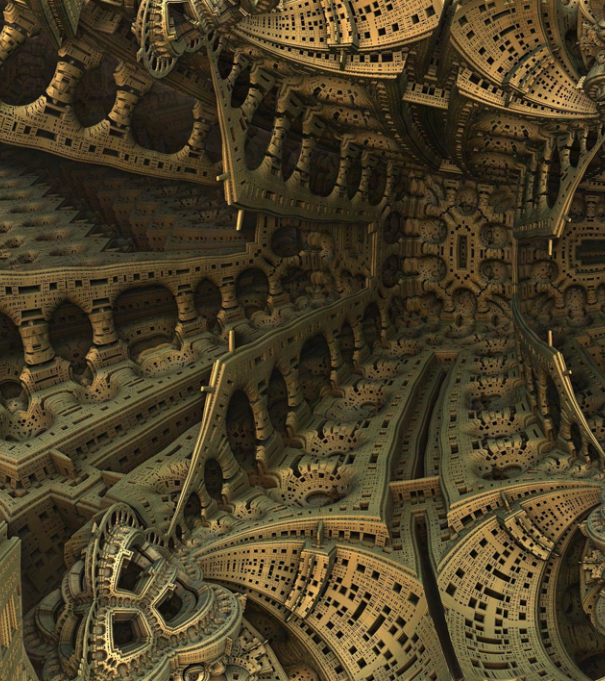
Overture. A brief historical case.

Act I. On conjectures, refutations, and argumentation.

Act II. There is no certain datum in the world.

Act III. Interesting problems are complex.

Epilogue.



A Trustworthy Loss Function

Classification becomes regression outputting pieces of evidences in favour of different classes

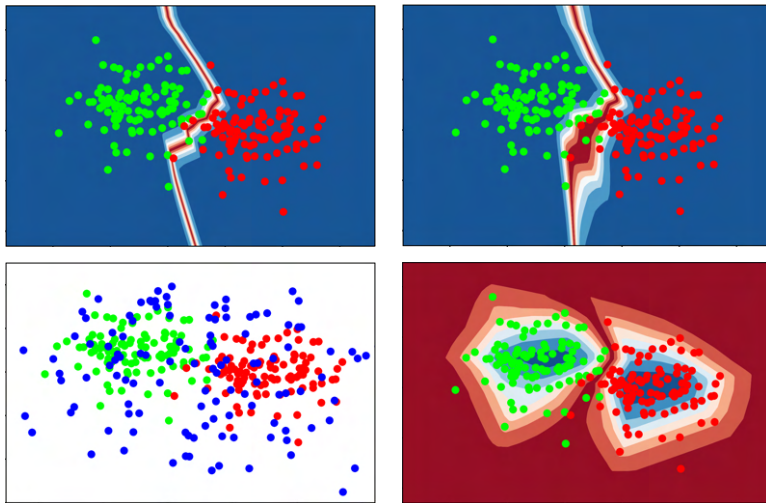
Expected squared error (aka Brier score) with $\text{Dir}(\mathbf{m}_i \mid \boldsymbol{\alpha}_i)$ (prior for a Multinomial) penalising the divergence from the uniform distribution:

$$\mathcal{L} = \sum_{i=1}^N \mathbb{E}[\|\mathbf{y}_i - \mathbf{m}_i\|_2^2] + \lambda_t \sum_{i=1}^N \text{KL}(\text{Dir}(\boldsymbol{\mu}_i \mid \tilde{\boldsymbol{\alpha}}_i) \parallel \text{Dir}(\boldsymbol{\mu}_i \mid \mathbf{1}))$$

where:

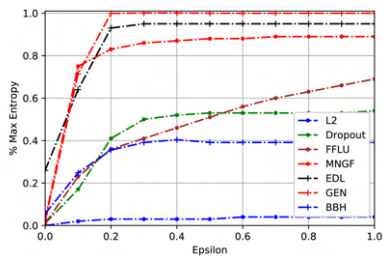
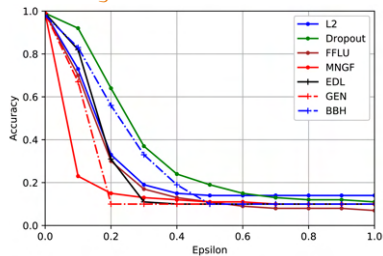
- λ_t avoid premature convergence to the uniform distribution;
- $\tilde{\boldsymbol{\alpha}}_i = \mathbf{y}_i + (\mathbf{1} - \mathbf{y}_i) \cdot \boldsymbol{\alpha}_i$ are the Dirichlet parameters the neural network in a forward pass has put on the wrong classes, and the idea is to minimise them as much as possible.
- $\text{KL}(\text{Dir}(\boldsymbol{\mu}_i \mid \tilde{\boldsymbol{\alpha}}_i) \parallel \text{Dir}(\boldsymbol{\mu}_i \mid \mathbf{1})) = \ln \left(\frac{\Gamma(\sum_{k=1}^K \tilde{\alpha}_{i,k})}{\Gamma(K) \prod_{k=1}^K \Gamma(\tilde{\alpha}_{i,k})} \right) + \sum_{k=1}^K (\tilde{\alpha}_{i,k} - 1) \left[\psi(\tilde{\alpha}_{i,k}) - \psi \left(\sum_{j=1}^K \tilde{\alpha}_{i,j} \right) \right]$
where $\psi(x) = \frac{d}{dx} \ln(\Gamma(x))$ is the *digamma* function

EDL + GAN for adversarial training

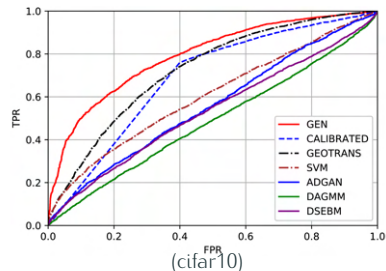
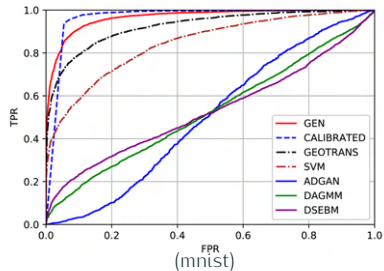


Şensoy, Kaplan, Cerutti, and Saleki. "Uncertainty-Aware Deep Classifiers using Generative Models." AAAI 2020

Robustness against FGS



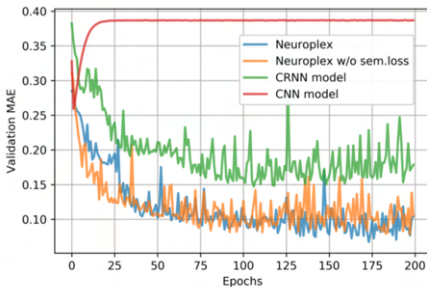
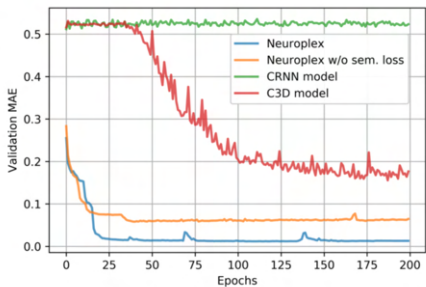
Anomaly detection



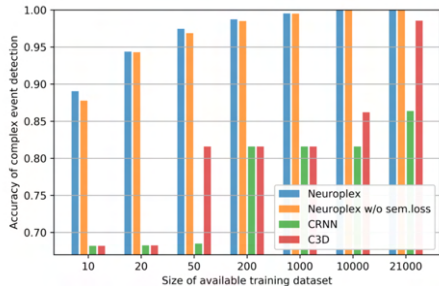
Şensoy, Kaplan, Cerutti, and Saleki. "Uncertainty-Aware Deep Classifiers using Generative Models." AAAI 2020

Roig Vilamala et. al. "A Hybrid Neuro-Symbolic Approach for Complex Event Processing (Extended Abstract)."
In ICLP2020.

Xing et. al. "Neuroplex: Learning to Detect Complex Events in Sensor Networks through Knowledge Injection."
In SenSys2020.



	Sim. 1	Sim. 2	Sim. 3	Sim. 4	Sim. 5
Window Length	10	20	30	3	2
# of Uniq Events	10	10	10	3	3
# of CE	4	4	7	5	4
Avg. CE Length	2.8	2.8	3.43	2	2
Neuroplex	99.39%	99.56%	98.65%	100.00%	99.98%
Lenet(Neuroplex)	98.87%	99.17%	98.91%	99.84%	99.78%
CRNN model	69.98%	7.79%	1.83%	86.37%	99.99%
C3D model	88.47%	83.73%	86.91%	98.56%	99.72%



Xing et. al. "Neuroplex: Learning to Detect Complex Events in Sensor Networks through Knowledge Injection."
In SenSys2020.

Overture. A brief historical case.

Act I. On conjectures, refutations, and argumentation.

Act II. There is no certain datum in the world.

Act III. Interesting problems are complex.

Epilogue.





www.grundnerco.at

www.grundnerco.at

65

70

85

90

95

100

105

110

115

120

125

130

135

140

145

150

155

160

165

170

175

180

185

190

Roig Vilamala et. al. "A Hybrid Neuro-Symbolic Approach for Complex Event Processing (Extended Abstract)."
In ICLP2020.

Xing et. al. "Neuroplex: Learning to Detect Complex Events in Sensor Networks through Knowledge Injection."
In SenSys2020.



Co-I

S. Chakraborty **IBM Research T. J. Watson** • M. Giacomini **Brescia** • L. Kaplan **US CDC ARL**
A. Kimmig **KU Leuven** • S. Julier **UCL** • Y. McDermott-Rees **Swansea** • T. Norman **Southampton**
N. Oren **Aberdeen** • G. Pearson **UK MoD Dstl** • A. Preece **Cardiff** • M. Şensoy **Ozyegin**
M. Srivastava **UCLA** • M. Thimm **Hagen** • N. Tintarev **Maastricht** • A. Toniolo **St. Andrews**
M. Vallati **Huddersfield**

Intern/PhD/Post-Doc

C. Allen **Cardiff** • A. Fanelli **Brescia** • L. Garcia **UCLA** • S. Habib **UCL** • C. Hougen **Michigan**
O. Lipinski **Southampton** • K. Mishra **US CDC ARL** • M. Roig Vilamala **Cardiff** • H. Rose **UCL**
G. Pellier-Hollows **Cardiff** • T. Xing **UCLA** • T. Zanetti **Cardiff**