

CAMS - EHESS

Séminaire Systèmes complexes en sciences sociales

2 Juin 2026 – Paris, France

Explicabilité formelle et interprétabilité théorique des modèles neuronaux

Juan Luis (Gianni) Gastaldi

ETH zürich

www.giannigastaldi.com

Intro: IA et sciences humaines

Explicabilité formelle

Interprétabilité théorique

Conclusion

Intro: IA et sciences humaines

Explicabilité formelle

Interprétabilité théorique

Conclusion

Behavior Revealed in Mobile Phone Usage Predicts Credit Repayment

Daniel Björkegren and Darrell Grissen

On Inferring User Socioeconomic Status with Mobility Records

Zheng Wang[†], Mingrui Liu^{*}, Cheng Long^{*¶}, Qianru Zhang[‡], Jiangneng Li^{*}, Chunyan Miao^{*§}

^{*}School of Computer Science and Engineering, Nanyang Technological University, Singapore,

[¶]Huawei Singapore Research Center, Singapore,

[‡]Department of Computer Science, The University of Hong Kong, Hong Kong SAR

[§]China-Singapore International Joint Research Institute (CSIJRI), China

Towards Deep Learning Models for Psychological State Prediction using Smartphone Data: Challenges and Opportunities

Gatis Mikelsons

University of Oxford and
The Alan Turing Institute

gatis.mikelsons@jesus.ox.ac.uk

Matthew Smith

University of Warwick and
The Alan Turing Institute

m.d.smith@warwick.ac.uk

Behavior Revealed in Mobile Phone Usage Predicts Credit Repayment

Daniel Björkegren and Darrell Grissen

On Inferring User Socioeconomic Status with Mobility Records

Zheng Wang⁺¹, Mingrui Liu^{*}, Cheng Long^{*4}, Qianru Zhang[‡], Jiangneng Li^{*}, Chunyan Miao^{*5}

^{*}School of Computer Science and Engineering, Nanyang Technological University, Singapore,

[‡]Huawei Singapore Research Center, Singapore,

[‡]Department of Computer Science, The University of Hong Kong, Hong Kong SAR

⁵China-Singapore International Joint Research Institute (CSIJRI), China

Towards Deep Learning Models for Psychological State Prediction using Smartphone Data: Challenges and Opportunities

Gatis Mikelsons
University of Oxford and
The Alan Turing Institute
gatis.mikelsons@jesus.ox.ac.uk

Matthew Smith
University of Warwick and
The Alan Turing Institute
m.d.smith@warwick.ac.uk

The Potential and Challenges of Evaluating Attitudes, Opinions, and Values in Large Language Models

Bolei Ma^{*LMU, m} Xinpeng Wang^{LMU, m} Tiancheng Hu[‡] Anna-Carolina Haensch^{LMU, 6}
Michael A. Hedderich^{LMU, m} Barbara Plank^{LMU, m, 6} Frauke Kreuter^{LMU, m, 6}

^{LMU}LMU Munich ^mMunich Center for Machine Learning [‡]University of Cambridge
⁶University of Maryland, College Park ⁶ITU Copenhagen

Multilingual Political Views of Large Language Models: Identification and Steering

Daniil Gurgurov^{1,5} Katharina Trinley¹ Ivan Vykopal^{3,4}
Josef van Genabith^{1,5} Simon Ostermann^{1,5} Roberto Zamparelli²

¹Saarland University ²University of Trento

³Brno University of Technology ⁴Kempelen Institute of Intelligent Technologies

⁵German Research Center for AI (DFKI)

The political ideology of conversational AI: Converging evidence on ChatGPT's pro-environmental, left-libertarian orientation

Jochen Hartmann^{a,1,2,3}, Jasper Schwenzow^{b,1}, and Maximilian Witte^{b,1}

^aTechnical University of Munich, TUM School of Management, Arcisstr. 21, 80333 Munich, Germany
^bUniversity of Hamburg, Hamburg Business School, Moorweidenstrasse 18, 20148 Hamburg, Germany

¹All authors contributed equally to this work.

Using sequences of life-events to predict human lives

Received: 6 June 2023

Accepted: 15 November 2023

Published online: 18 December 2023

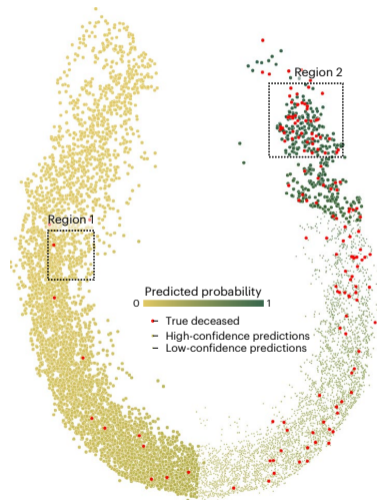
 Check for updates

Germans Savcisen¹, Tina Eliassi-Rad^{2,3}, Lars Kai Hansen¹,
Laust Hvas Mortensen^{4,5}, Lau Lilleholt^{6,7}, Anna Rogers⁸, Ingo Zettler^{6,7} &
Sune Lehmann^{1,7} ✉

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

Exemple: life2vec

d Person embedding space
(projected with PaCMAP)



Méthodes d'interprétabilité mécanistique (Geiger et al., 2025)

Behavioral Methods

- ◇ Feature attribution
- ◇ Integrated gradients
- ◇ Effects of real-world concepts on models

Patching Activations with Interchange Interventions

- ◇ Interchange interventions
- ◇ Path patching
- ◇ Causal mediation analysis

Ablation-Based Analysis

- ◇ Concept erasure
- ◇ Sub-circuit analysis
- ◇ Causal scrubbing

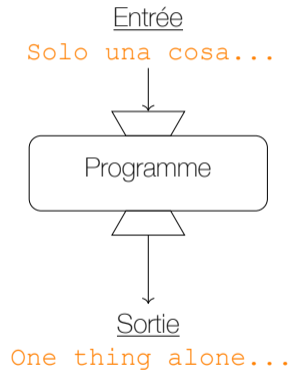
Modular Feature Learning

- ◇ Probing
- ◇ Sparse autoencoders
- ◇ PCA
- ◇ Differential Binary Masking
- ◇ Difference of means
- ◇ Distributed Alignment Search

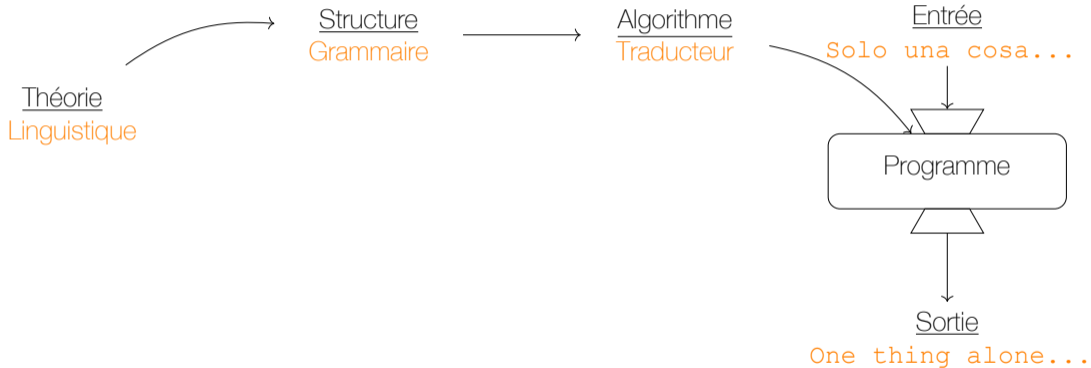
Other Approaches

- ◇ Activation Steering
- ◇ Training for Interpretability
- ◇ Causal Abstraction

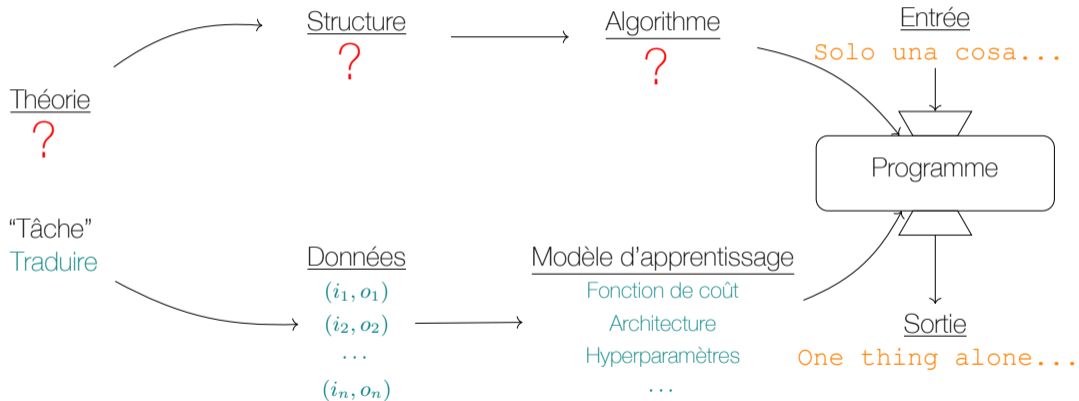
La structure implicite des données



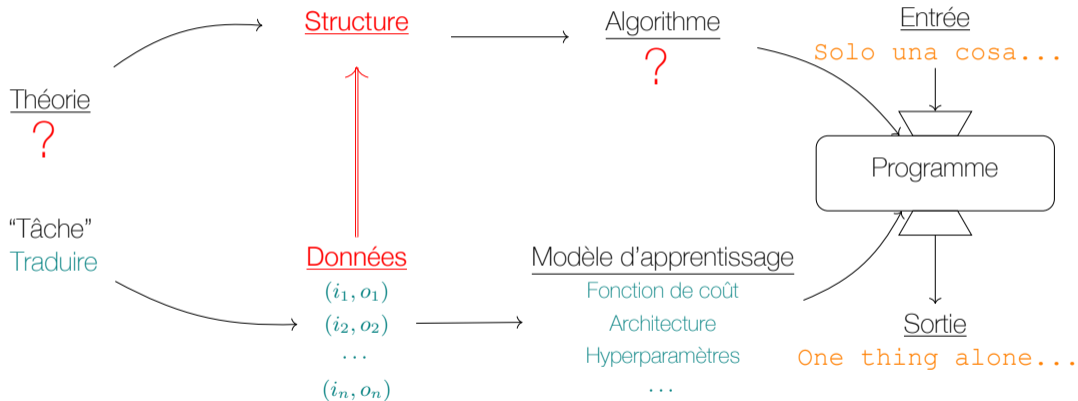
La structure implicite des données



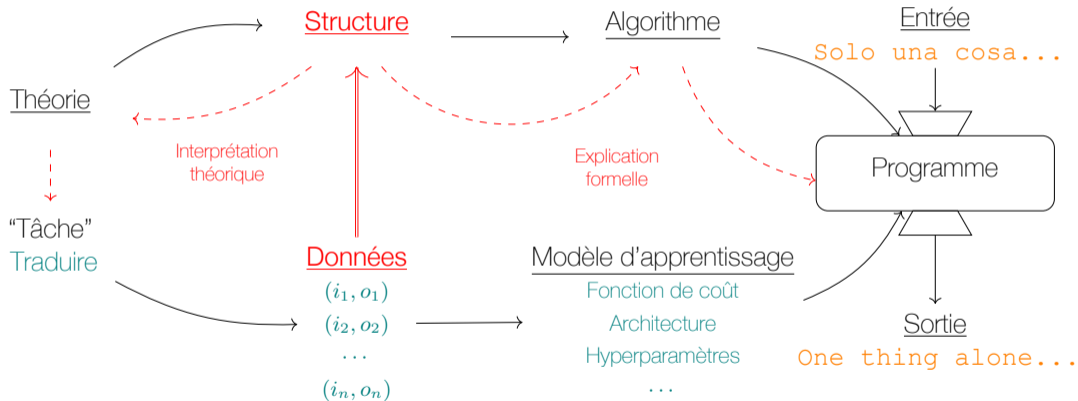
La structure implicite des données



La structure implicite des données



La structure implicite des données



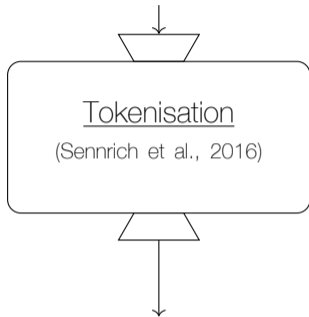
Intro: IA et sciences humaines

Explicabilité formelle

Interprétabilité théorique

Conclusion

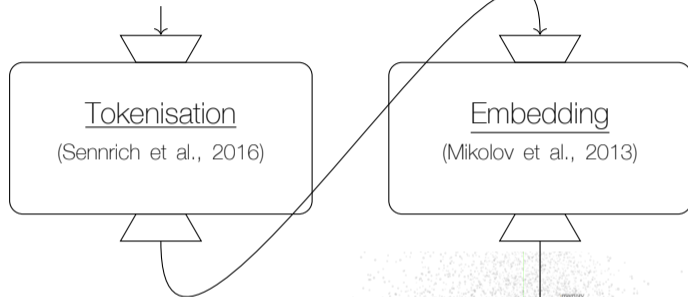
Epistemology of Machine Learning
Distributional Language Models



Epistemology of Machine Learning
Distributional Language Models

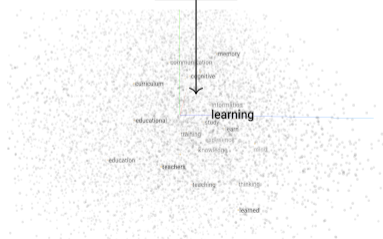
(<https://tiktokenizer.vercel.app>)

Epistemology of Machine Learning
Distributional Language Models



Epistemology of Machine Learning
Distributional Language Models

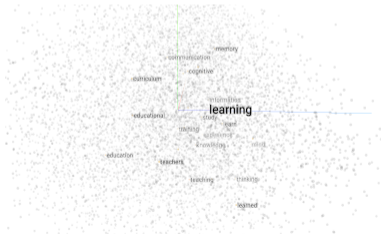
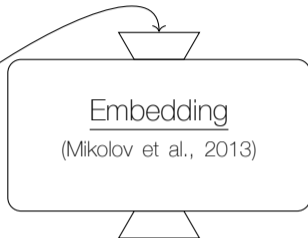
(<https://tiktokenizer.vercel.app>)



(<https://projector.tensorflow.org>)

Epistemology of Machine Learning
Distributional Language Models

(<https://tiktokenizer.vercel.app>)



(<https://projector.tensorflow.org>)

La structure des embeddings

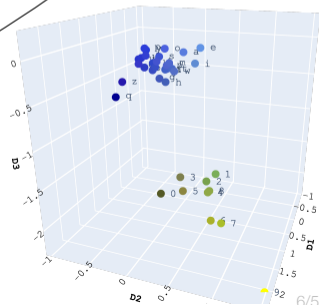
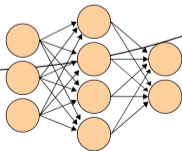
Structure

?

{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}

Embedding

Données



word2vec expliqué (Levy and Goldberg, 2014)

$$\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)]$$

Où:

\vec{w} = représentation vectorielle du mot w

\vec{c} = représentation vectorielle du contexte c

$\sigma(x)$ = $\frac{1}{1+e^{-x}}$

k = nombre d'échantillons "négatifs" (arbitraires)

c_N = contexte arbitraire tiré de P_D

$P_D(c)$ = distribution unigramme empirique de c dans les données D , c-à-d $\frac{\#(c)}{|D|}$

word2vec expliqué (Levy and Goldberg, 2014)

$$\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)]$$

word2vec expliqué (Levy and Goldberg, 2014)

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

word2vec expliqué (Levy and Goldberg, 2014)

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

$$\frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0$$

word2vec expliqué (Levy and Goldberg, 2014)

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

$$\frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{quand} \quad \vec{w} \cdot \vec{c} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

word2vec expliqué (Levy and Goldberg, 2014)

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

$$\begin{aligned} \frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{quand} \quad \vec{w} \cdot \vec{c} &= \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k \\ &= \text{PMI}(w, c) - \log k \end{aligned}$$

word2vec expliqué (Levy and Goldberg, 2014)

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

$$\begin{aligned} \frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{quand} \quad \vec{w} \cdot \vec{c} &= \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k \\ &= \text{PMI}(w, c) - \log k \end{aligned}$$

Contrainte supplémentaire :

\vec{w} et \vec{c} doivent être **de faible dimension**

word2vec expliqué (Levy and Goldberg, 2014)

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

$$\begin{aligned} \frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{quand} \quad \vec{w} \cdot \vec{c} &= \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k \\ &= \text{PMI}(w, c) - \log k \end{aligned}$$

Contrainte supplémentaire :

\vec{w} et \vec{c} doivent être **de faible dimension**

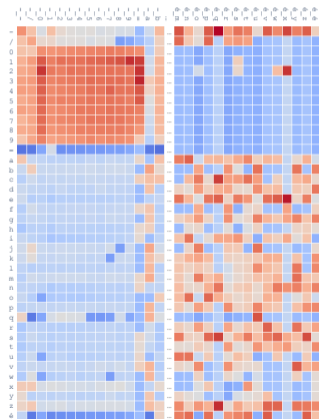
La **décomposition en valeurs singulières (SVD)** fournit une **solution exacte** à ce problème d'optimisation.

Exemple: Caractères dans Wikipédia

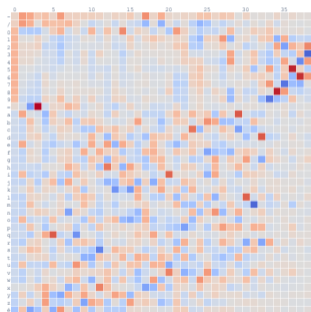
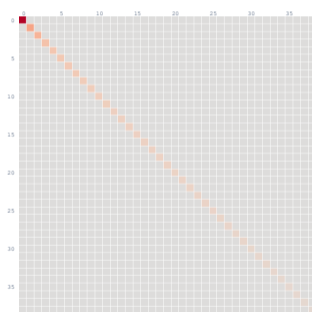
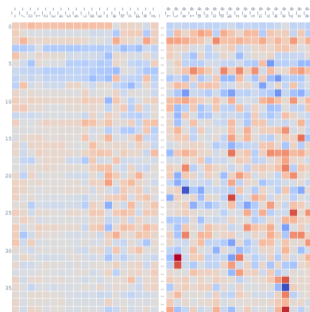
$W = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$

$C = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$

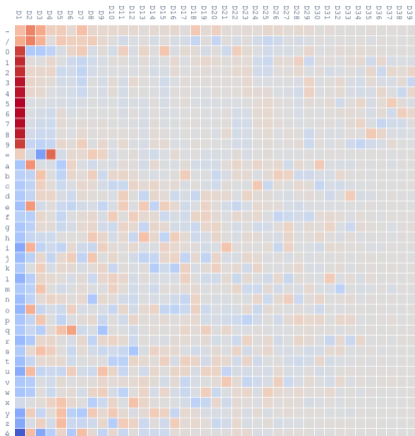
$$\begin{aligned}M_{wc} &= \text{pmi}(w, c) \\ &= \log \frac{p(w, c)}{p(w)p(c)}\end{aligned}$$



SVD d'une matrice pmi des caractères dans Wikipédia

 U  Σ  V^T 

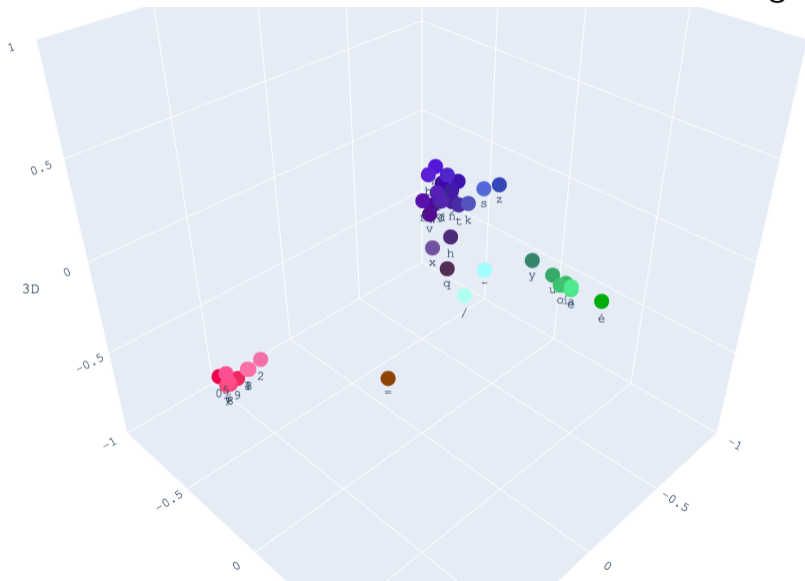
$$U \times \Sigma$$



$$\hat{U} \times \hat{\Sigma}$$



$$\hat{U} \times \hat{\Sigma}$$



La structure des embeddings

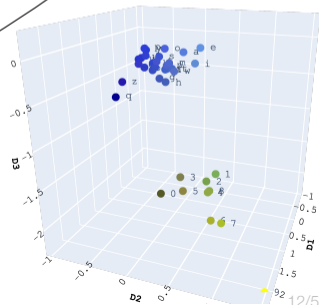
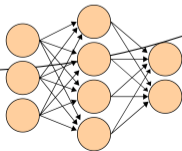
Structure

?

{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}

Embedding

Données

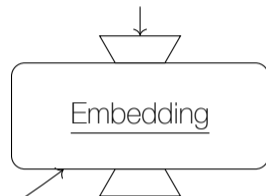


La structure des embeddings

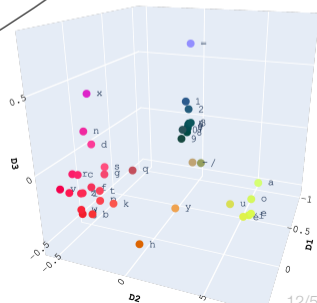
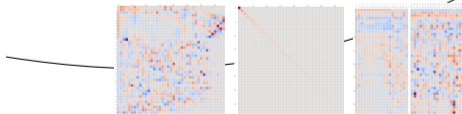
Structure



{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}



Données



4 Why does this produce good word representations?

Good question. We don't really know.

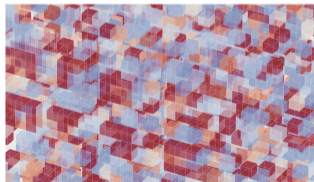
The distributional hypothesis states that words in similar contexts have similar meanings. The objective above clearly tries to increase the quantity $v_w \cdot v_c$ for good word-context pairs, and decrease it for bad ones. Intuitively, this means that words that share many contexts will be similar to each other (note also that contexts sharing many words will also be similar to each other). This is, however, very hand-wavy.

Can we make this intuition more precise? We'd really like to see something more formal.

(Goldberg and Levy, 2014)

The Structure of Meaning in Language: Parallel Narratives in Linear Algebra and Category Theory

*Tai-Danae Bradley, Juan Luis Gastaldi,
and John Terilla*



Introduction

Categories for AI, an online program about category the-

intelligence in particular. While this article is by no means a comprehensive report on that event, the popularity of “Cats for AI” — the five introductory lectures have been viewed thousands of times — signals the growing prevalence of category theoretic tools in AI.

One way that category theory is gaining traction in machine learning is by providing a formal way to discuss how learning systems can be put together. This article has a different and somewhat narrow focus. It’s about how a fundamental piece of AI technology used in language modeling can be understood, with the aid of categorical thinking, as a process that extracts structural features of language from purely syntactical input. The idea that structure arises from form may not be a surprise for many readers — category theoretic ideas have been a major influence in pure

Vecteurs de mots comme fonctions sur des ensembles

$$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$$

$$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$$

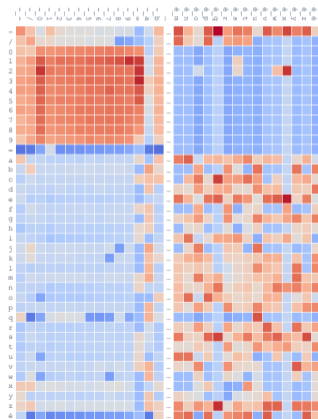
Vecteurs de mots comme fonctions sur des ensembles

$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$

$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$

$M: X \times Y \rightarrow \mathbb{R}$

$(x, y) \mapsto \text{pmi}(x, y)$



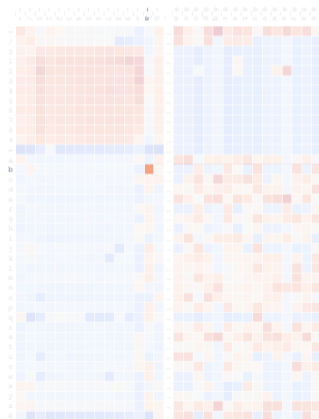
Vecteurs de mots comme fonctions sur des ensembles

$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$

$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$

$M: X \times Y \rightarrow \mathbb{R}$

$(x, y) \mapsto \text{pmi}(x, y)$



Vecteurs de mots comme fonctions sur des ensembles

$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$

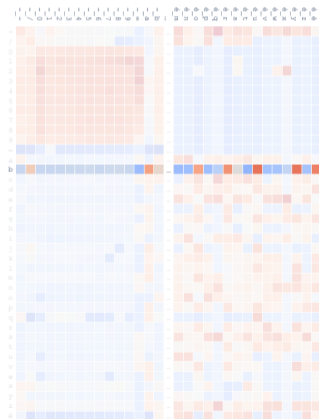
$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$

$M: X \times Y \rightarrow \mathbb{R}$

$(x, y) \mapsto \text{pmi}(x, y)$

$M_x: X \rightarrow \mathbb{R}^Y$

$x \mapsto M(x, -)$



Vecteurs de mots comme fonctions sur des ensembles

$$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$$

$$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$$

$$M: X \times Y \rightarrow \mathbb{R}$$

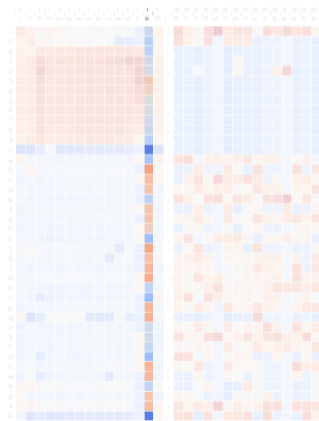
$$(x, y) \mapsto \text{pmi}(x, y)$$

$$M_x: X \rightarrow \mathbb{R}^Y$$

$$x \mapsto M(x, -)$$

$$M_y: Y \rightarrow \mathbb{R}^X$$

$$y \mapsto M(-, y)$$



Vecteurs de mots comme fonctions sur des ensembles

$$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$$

$$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$$

$$M: X \times Y \rightarrow \mathbb{R}$$

$$(x, y) \mapsto \text{pmi}(x, y)$$

$$X \xrightarrow{M_x} \mathbb{R}^Y$$

$$M_x: X \rightarrow \mathbb{R}^Y$$

$$x \mapsto M(x, -)$$

$$\mathbb{R}^X \xleftarrow{M_y} Y$$

$$M_y: Y \rightarrow \mathbb{R}^X$$

$$y \mapsto M(-, y)$$

Vecteurs de mots comme fonctions sur des ensembles

$$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$$

$$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$$

$$M: X \times Y \rightarrow \mathbb{R}$$

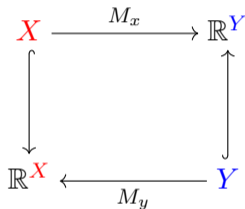
$$(x, y) \mapsto \text{pmi}(x, y)$$

$$M_x: X \rightarrow \mathbb{R}^Y$$

$$x \mapsto M(x, -)$$

$$M_y: Y \rightarrow \mathbb{R}^X$$

$$y \mapsto M(-, y)$$


$$\begin{array}{ccc} X & \xrightarrow{M_x} & \mathbb{R}^Y \\ \downarrow & & \uparrow \\ \mathbb{R}^X & \xleftarrow{M_y} & Y \end{array}$$

Vecteurs de mots comme fonctions sur des ensembles

$$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, \acute{e}\}$$

$$Y = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\acute{e}, z), (\acute{e}, \acute{e})\}$$

$$M: X \times Y \rightarrow \mathbb{R}$$

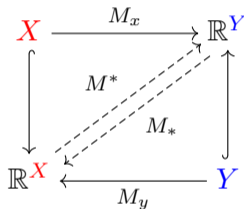
$$(x, y) \mapsto \text{pmi}(x, y)$$

$$M_x: X \rightarrow \mathbb{R}^Y$$

$$x \mapsto M(x, -)$$

$$M_y: Y \rightarrow \mathbb{R}^X$$

$$y \mapsto M(-, y)$$

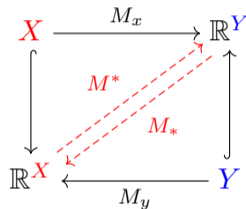


$$M^*: \mathbb{R}^X \rightarrow \mathbb{R}^Y$$

$$M_*: \mathbb{R}^Y \rightarrow \mathbb{R}^X$$

Vecteurs de mots comme fonctions sur des ensembles

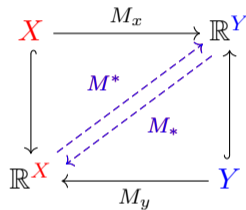
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$



Vecteurs de mots comme fonctions sur des ensembles

$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$



Vecteurs de mots comme fonctions sur des ensembles

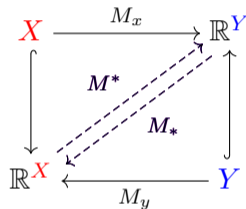
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$

$$\{u_1, \dots, u_m\} \subset \mathbb{R}^X$$

$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

$$\{\lambda_1, \dots, \lambda_{\min(m,n)}, 0, \dots, 0\}$$



Vecteurs de mots comme fonctions sur des ensembles

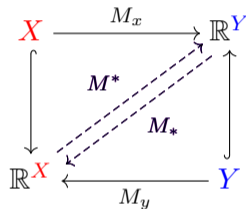
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$

$$\{u_1, \dots, u_m\} \subset \mathbb{R}^X$$

$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

$$\{\lambda_1, \dots, \lambda_{\min(m,n)}, 0, \dots, 0\}$$



$$U := [u_1, \dots, u_m]$$

$$M = U \Sigma V^T \quad V := [v_1, \dots, v_n]$$

$$\Sigma := \begin{bmatrix} \sqrt{\lambda_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sqrt{\lambda_r} \end{bmatrix}$$

Vecteurs de mots comme fonctions sur des ensembles

$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$

$$\{u_1, \dots, u_m\} \subset \mathbb{R}^X$$

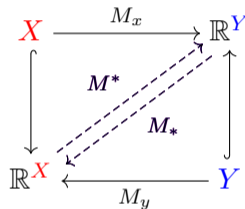
$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

$$\{\lambda_1, \dots, \lambda_{\min(m,n)}, 0, \dots, 0\}$$

$$M_* M^* u_i = \lambda_i u_i$$

$$M^* M_* v_i = \lambda_i v_i$$

Les u_i and v_i sont des **points fixes** (linéaires)!



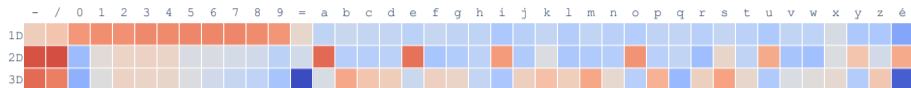
$$M = U \Sigma V^T$$

$$U := [u_1, \dots, u_m]$$

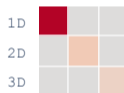
$$V := [v_1, \dots, v_n]$$

$$\Sigma := \begin{bmatrix} \sqrt{\lambda_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sqrt{\lambda_r} \end{bmatrix}$$

Vecteurs propres de M_*M^* :

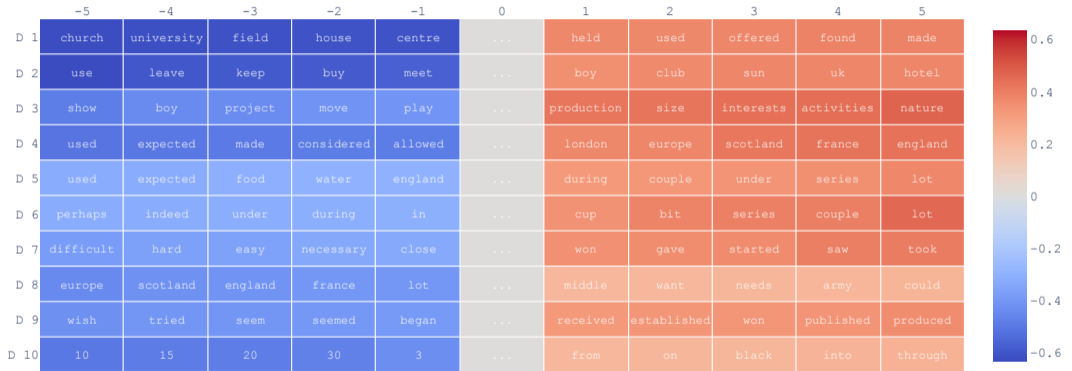


Valeurs propres de M_*M^* et M^*M_* :



Vecteurs propres de M^*M_* :





La structure des embeddings

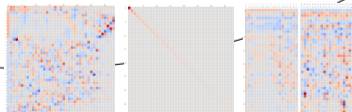
Structure

?

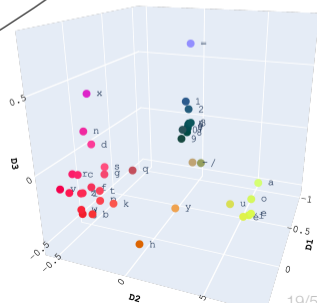
{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}

Embedding

Données



SVD

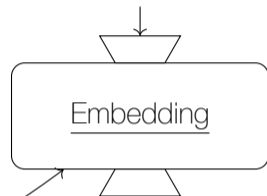


La structure des embeddings

Structure

?

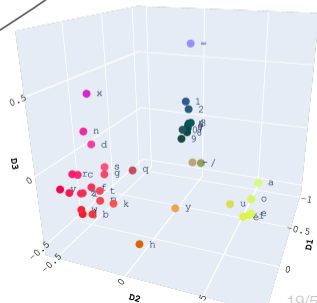
{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}



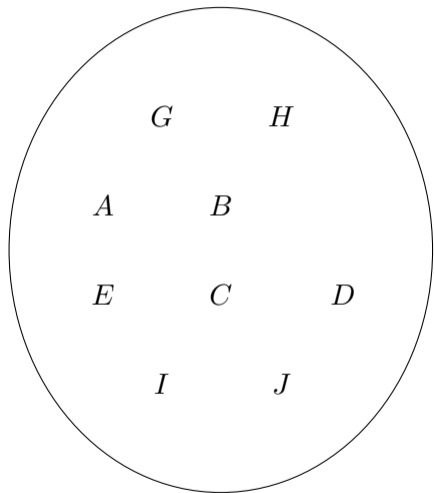
Données



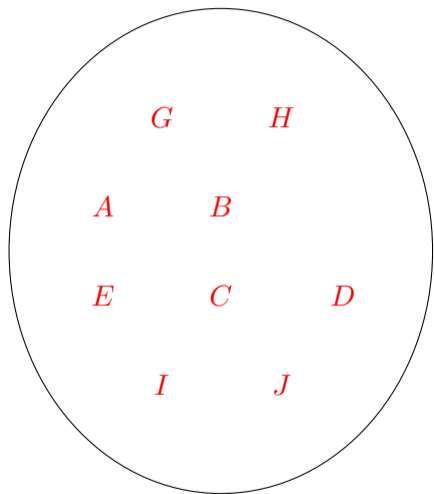
$$C^{\text{op}} \times D \rightarrow \mathbb{R}^i$$



Une catégorie est comme un ensemble muni d'une structure



Une catégorie est comme un ensemble muni d'une structure

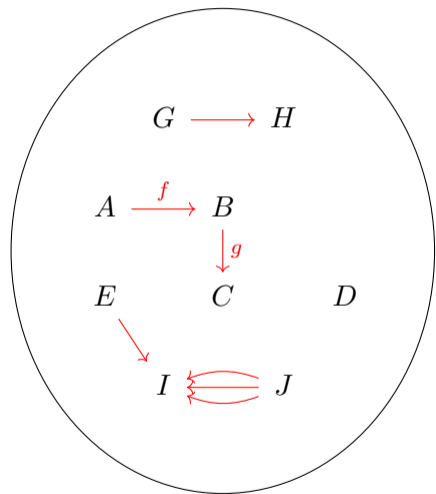


Definition (Category – Awodey, 2010)

Data:

- ◇ Objects: A, B, C, \dots

Une catégorie est comme un ensemble muni d'une structure

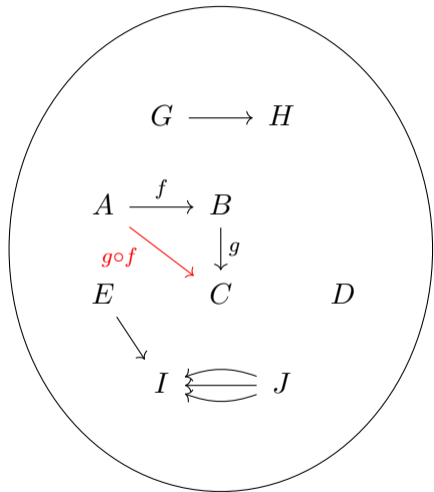


Definition (Category – Awodey, 2010)

Data:

- ◇ Objects: A, B, C, \dots
- ◇ Arrows: f, g, \dots

Une catégorie est comme un ensemble muni d'une structure



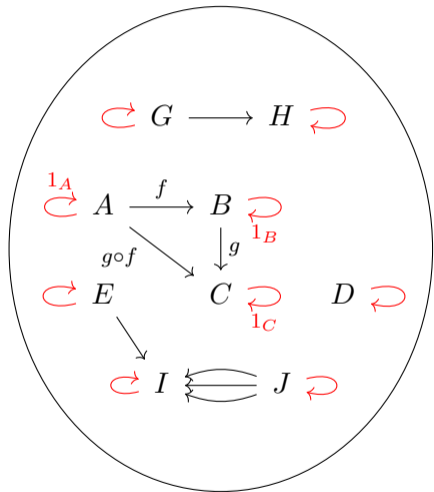
Definition (Category – Awodey, 2010)

Data:

- ◇ Objects: A, B, C, \dots
- ◇ Arrows: f, g, \dots
- ◇ Composition: Given $f : A \rightarrow B$ and $g : B \rightarrow C$, there is given an arrow

$$g \circ f : A \rightarrow C$$

Une catégorie est comme un ensemble muni d'une structure



Definition (Category – Awodey, 2010)

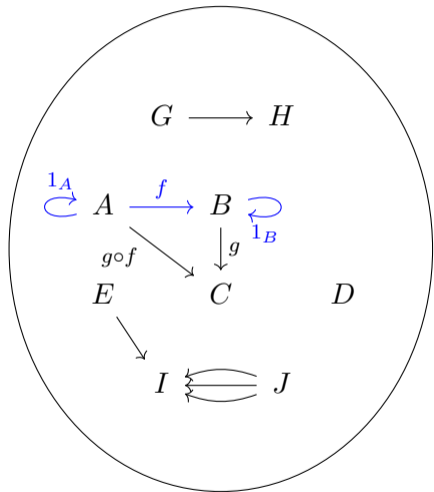
Data:

- ◇ Objects: A, B, C, \dots
- ◇ Arrows: f, g, \dots
- ◇ Composition: Given $f : A \rightarrow B$ and $g : B \rightarrow C$, there is given an arrow

$$g \circ f : A \rightarrow C$$

- ◇ Identity: For each A , there is $1_A : A \rightarrow A$

Une catégorie est comme un ensemble muni d'une structure



Definition (Category – Awodey, 2010)

Data:

- ◇ Objects: A, B, C, \dots
- ◇ Arrows: f, g, \dots
- ◇ Composition: Given $f : A \rightarrow B$ and $g : B \rightarrow C$, there is given an arrow

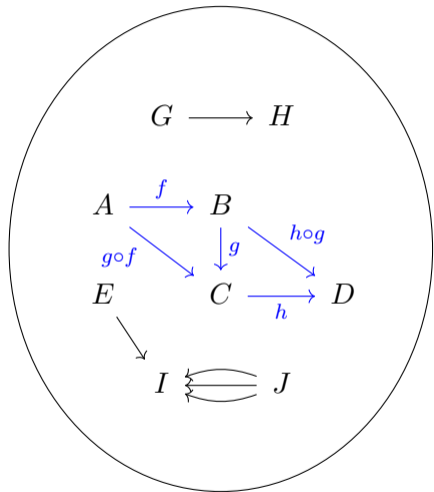
$$g \circ f : A \rightarrow C$$

- ◇ Identity: For each A , there is $1_A : A \rightarrow A$

Laws:

- ◇ Unit: $f \circ 1_A = f = 1_B \circ f$

Une catégorie est comme un ensemble muni d'une structure



Definition (Category – Awodey, 2010)

Data:

- ◇ Objects: A, B, C, \dots
- ◇ Arrows: f, g, \dots
- ◇ Composition: Given $f : A \rightarrow B$ and $g : B \rightarrow C$, there is given an arrow

$$g \circ f : A \rightarrow C$$

- ◇ Identity: For each A , there is $1_A : A \rightarrow A$

Laws:

- ◇ Unit: $f \circ 1_A = f = 1_B \circ f$
- ◇ Associativity: $f \circ (g \circ h) = (f \circ g) \circ h$

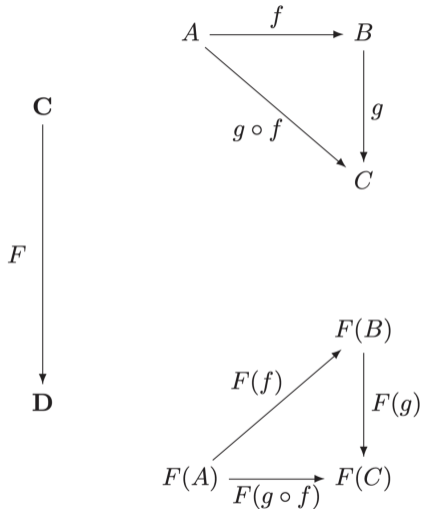
Un foncteur est une application entre catégories

Definition (Functor – Awodey, 2010)

A functor

$$F: \mathbf{C} \rightarrow \mathbf{D}$$

between categories \mathbf{C} and \mathbf{D} is a mapping of objects to objects and arrows to arrows, in such a way that



Un foncteur est une application entre catégories

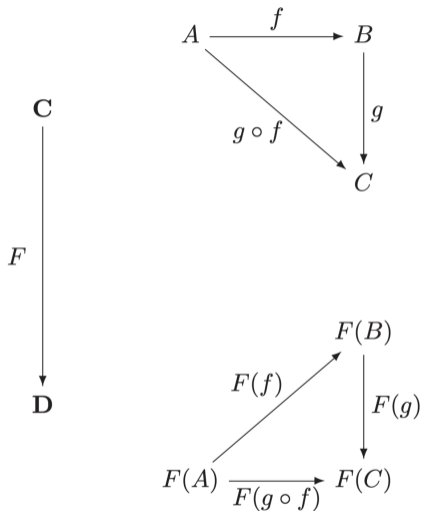
Definition (Functor – Awodey, 2010)

A functor

$$F: \mathbf{C} \rightarrow \mathbf{D}$$

between categories \mathbf{C} and \mathbf{D} is a mapping of objects to objects and arrows to arrows, in such a way that

(a) $F(f : A \rightarrow B) = F(f) : F(A) \rightarrow F(B)$



Un foncteur est une application entre catégories

Definition (Functor – Awodey, 2010)

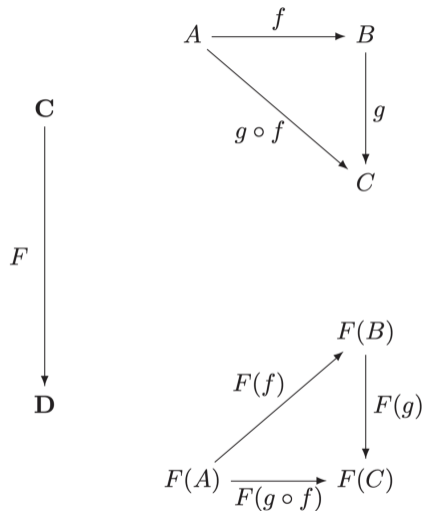
A functor

$$F: \mathbf{C} \rightarrow \mathbf{D}$$

between categories \mathbf{C} and \mathbf{D} is a mapping of objects to objects and arrows to arrows, in such a way that

(a) $F(f : A \rightarrow B) = F(f) : F(A) \rightarrow F(B)$

(b) $F(1_A) = 1_{F(A)}$



Un foncteur est une application entre catégories

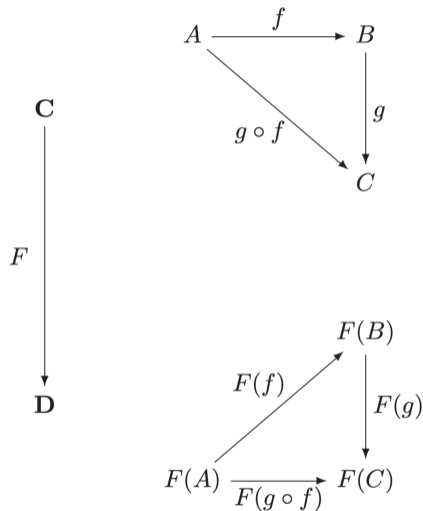
Definition (Functor – Awodey, 2010)

A functor

$$F: \mathbf{C} \rightarrow \mathbf{D}$$

between categories \mathbf{C} and \mathbf{D} is a mapping of objects to objects and arrows to arrows, in such a way that

- (a) $F(f : A \rightarrow B) = F(f) : F(A) \rightarrow F(B)$
- (b) $F(1_A) = 1_{F(A)}$
- (c) $F(g \circ f) = F(g) \circ F(f)$



Definition 2.15. In any category \mathbf{C} , a *product diagram* for the objects A and B consists of an object P and arrows

$$A \xleftarrow{p_1} P \xrightarrow{p_2} B$$

satisfying the following UMP:

Given any diagram of the form

$$A \xleftarrow{x_1} X \xrightarrow{x_2} B$$

there exists a unique $u : X \rightarrow P$, making the diagram

$$\begin{array}{ccccc} & & X & & \\ & \swarrow & \vdots & \searrow & \\ & x_1 & u & x_2 & \\ & \swarrow & \downarrow & \searrow & \\ A & \xleftarrow{p_1} & P & \xrightarrow{p_2} & B \end{array}$$

commute, that is, such that $x_1 = p_1 u$ and $x_2 = p_2 u$.

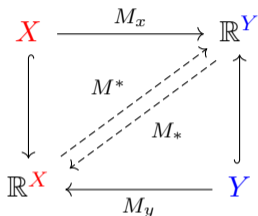
(Awodey, 2010)

Un profoncteur est un foncteur du produit de deux catégories quelconques vers la catégorie Set

$$\begin{array}{ccc} \textit{term}_i & \textit{context}_i & \textit{measure} \\ \downarrow & \downarrow & \swarrow \\ \mathbf{C}^{\text{op}} & \times \mathbf{D} & \rightarrow \mathbf{Set} \end{array}$$

Opérateur "distributionnel" et points fixes

$$M : X \times Y \rightarrow \mathbb{R}$$



$$\left. \begin{array}{l} M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X \\ M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y \end{array} \right\} \text{opérateur} \\ \text{distributionnel}$$

$$\left. \begin{array}{l} \{u_i \in \mathbb{R}^X \mid M_* M^* u_i = \lambda_i u_i\} \\ \{v_i \in \mathbb{R}^Y \mid M^* M_* v_i = \lambda_i v_i\} \end{array} \right\} \text{points} \\ \text{fixes}$$

Opérateur "distributionnel" et points fixes

$$M : X \times Y \rightarrow \mathbb{R}$$

$$\begin{array}{ccc}
 X & \xrightarrow{M_x} & \mathbb{R}^Y \\
 \downarrow & \nearrow M^* & \uparrow \\
 \mathbb{R}^X & & Y \\
 & \xleftarrow{M_y} &
 \end{array}$$

$$\left. \begin{array}{l}
 M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X \\
 M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y
 \end{array} \right\} \text{opérateur distributionnel}$$

$$\left. \begin{array}{l}
 \{u_i \in \mathbb{R}^X \mid M_* M^* u_i = \lambda_i u_i\} \\
 \{v_i \in \mathbb{R}^Y \mid M^* M_* v_i = \lambda_i v_i\}
 \end{array} \right\} \text{points fixes}$$

$$\mathbf{C}^{\text{op}} \times \mathbf{D} \rightarrow \text{Set}$$

$$\begin{array}{ccc}
 \mathbf{C} & \xrightarrow{\mathcal{M}_c} & (\text{Set}^{\mathbf{D}})^{\text{op}} \\
 \downarrow \text{Yoneda} & \nearrow \mathcal{M}^* & \uparrow \text{Yoneda} \\
 \text{Set}^{\mathbf{C}^{\text{op}}} & & \mathbf{D} \\
 & \xleftarrow{\mathcal{M}_d} &
 \end{array}$$

$$\left\{ \begin{array}{l}
 \mathcal{M}_* \mathcal{M}^* : \text{Set}^{\mathbf{C}^{\text{op}}} \rightarrow \text{Set}^{\mathbf{C}^{\text{op}}} \\
 \mathcal{M}^* \mathcal{M}_* : (\text{Set}^{\mathbf{D}})^{\text{op}} \rightarrow (\text{Set}^{\mathbf{D}})^{\text{op}}
 \end{array} \right.$$

$$\left\{ \begin{array}{l}
 \{f \in \text{Set}^{\mathbf{C}^{\text{op}}} \mid \mathcal{M}_* \mathcal{M}^*(f) \cong f\} \\
 \{g \in (\text{Set}^{\mathbf{D}})^{\text{op}} \mid \mathcal{M}^* \mathcal{M}_*(g) \cong g\}
 \end{array} \right.$$

Profoncteur de catégories enrichies

$$\begin{array}{ccc} \textit{term}_i & \textit{context}_i & \textit{measure} \\ \downarrow & \downarrow & \swarrow \\ \mathbf{C}^{\text{op}} & \times \mathbf{D} & \rightarrow \mathbf{Set} \end{array}$$

Profoncteur de catégories enrichies

$$\begin{array}{ccc} \textit{term}_i & \textit{context}_i & \textit{measure} \\ \downarrow & \downarrow & \swarrow \\ \mathbf{C}^{\text{op}} \times \mathbf{D} & \rightarrow & \mathbf{V} \end{array}$$

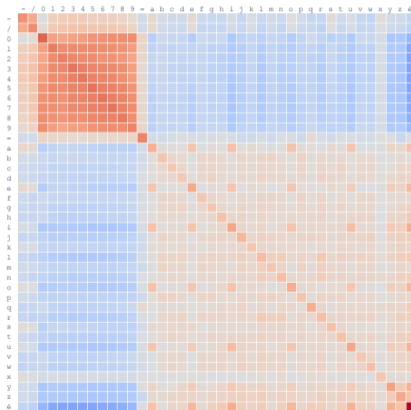
Profoncteur de catégories enrichies

$$\begin{array}{ccc} \textit{term}_i & \textit{context}_i & \textit{measure} \\ \downarrow & \downarrow & \swarrow \\ \mathbf{C}^{\text{op}} & \times \mathbf{D} & \rightarrow \mathbf{2} \end{array}$$

Profoncteur de catégories enrichies

$$\begin{array}{ccc} \text{term}_i & \text{context}_i & \text{measure} \\ \downarrow & \downarrow & \swarrow \\ \mathbf{C}^{\text{op}} \times \mathbf{D} & \rightarrow & \mathbf{2} \\ \Downarrow & & \\ \mathcal{M}^* : \mathbf{2}^{\mathbf{C}^{\text{op}}} & \rightleftarrows & (\mathbf{2}^{\mathbf{D}})^{\text{op}} : \mathcal{M}_* \end{array}$$

$$M_* M^* u = \lambda u$$



×

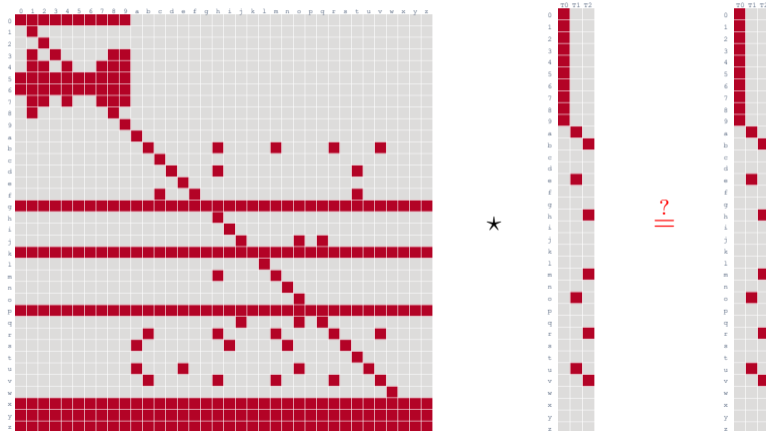


=

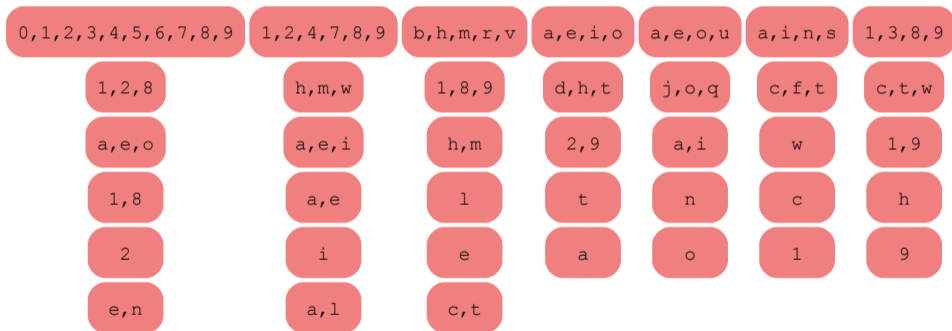


$$\begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{pmatrix}$$

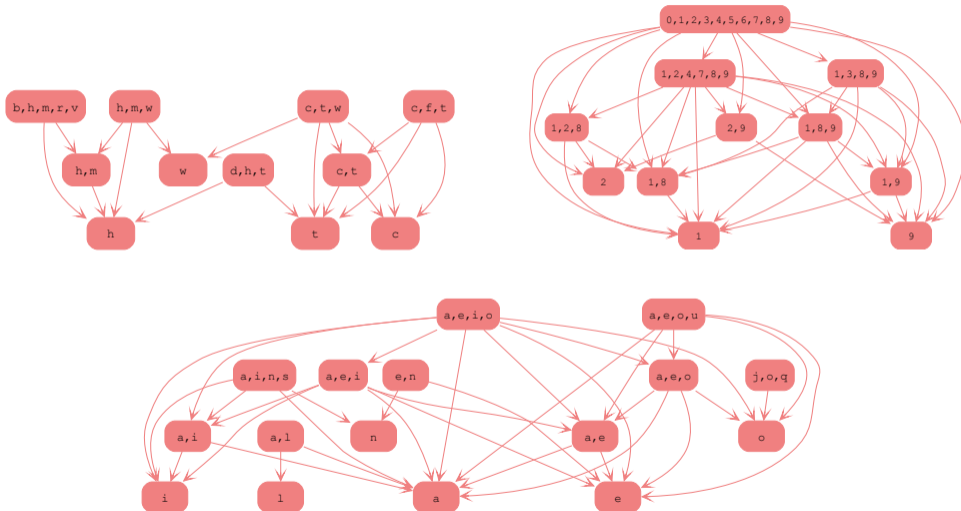
$$M_* M^* f = f$$



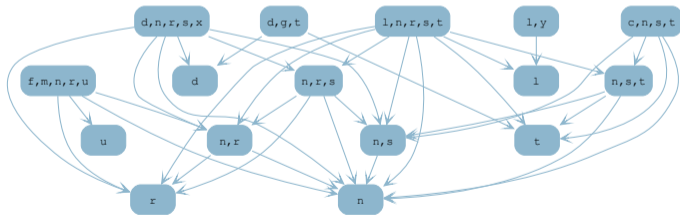
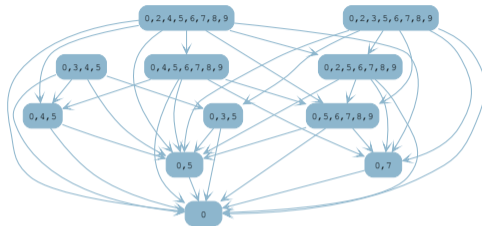
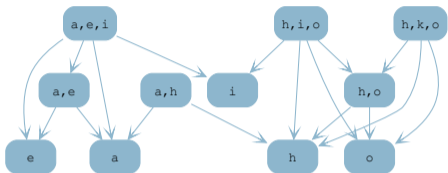
$$\mathcal{M}_* \mathcal{M}^* f = f$$



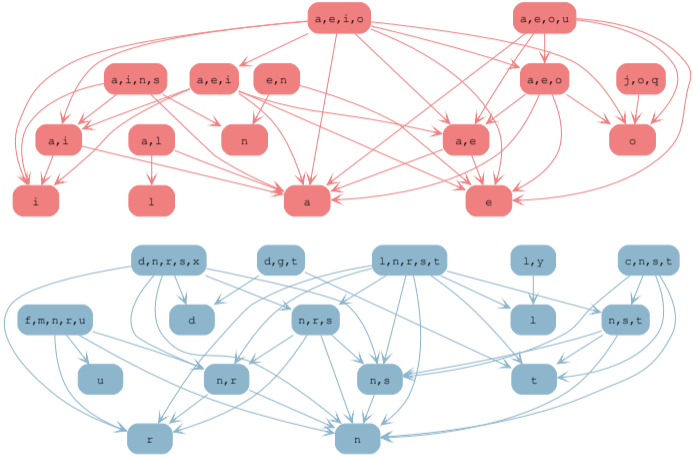
Structure d'ordre partiel



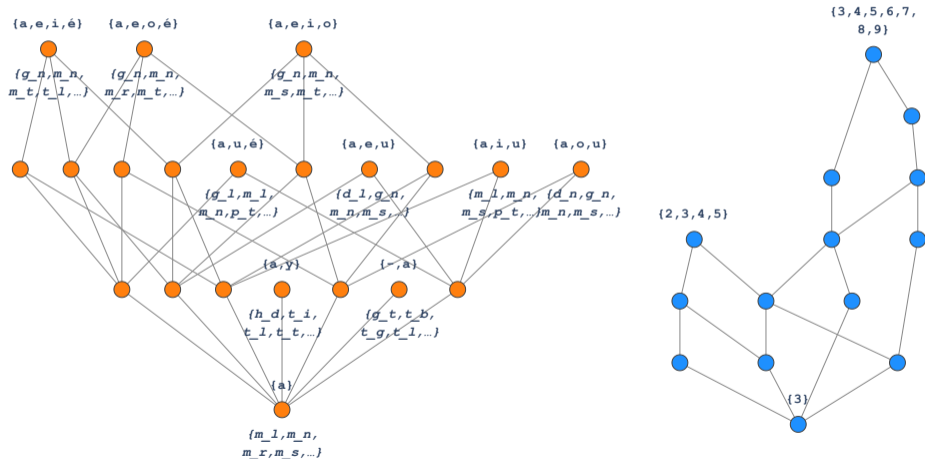
Ordre partiel dual



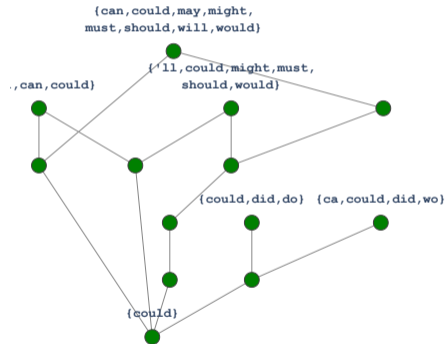
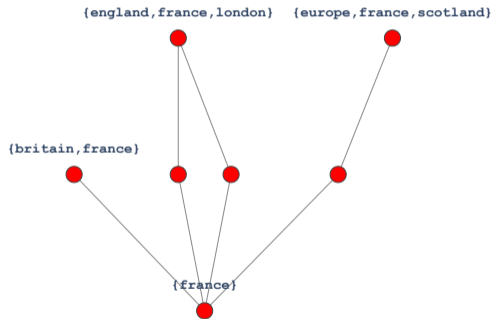
Couplage des points fixes



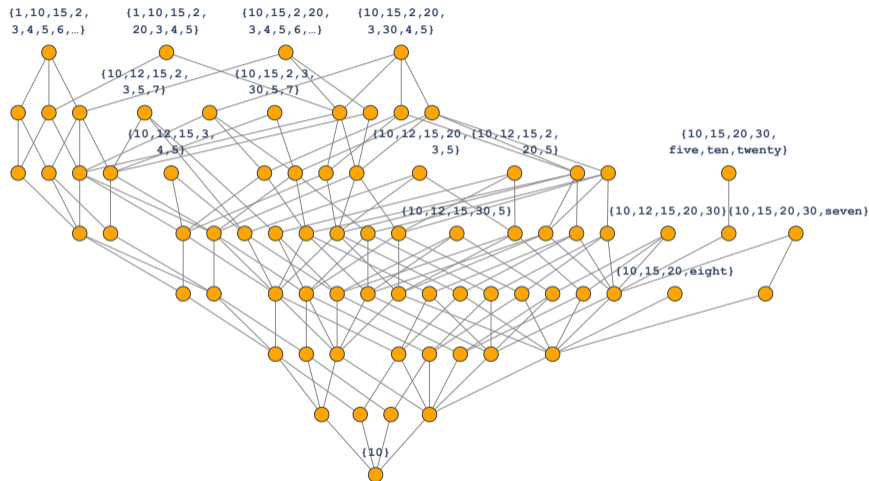
Concepts formels



Concepts formels (mots)



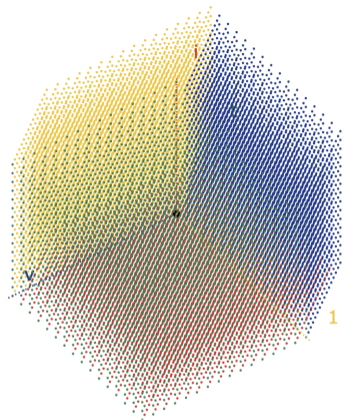
Concepts formels (mots)



Profoncteur et structures du noyau

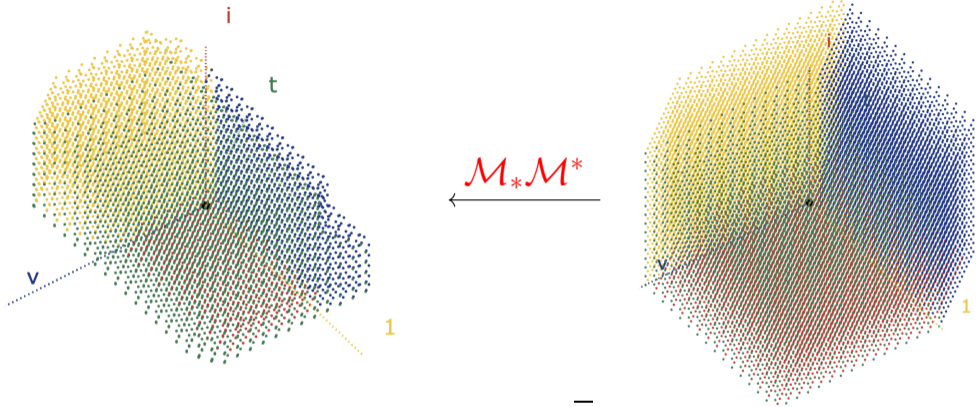
$$\begin{array}{ccc} e_i & s_i & \text{mesure} \\ \downarrow & \downarrow & \downarrow \\ \mathbb{C}^{\text{op}} \times \mathbb{D} & \rightarrow & \bar{\mathbb{R}} \\ \Downarrow & & \\ \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbb{C}^{\text{op}}} & \rightleftharpoons & (\bar{\mathbb{R}}^{\mathbb{D}})^{\text{op}} : \mathcal{M}_* \end{array}$$

Profoncteur et structures du noyau



$$\begin{array}{ccc} \mathbb{C}^{\text{op}} \times \mathbb{D} & \rightarrow & \bar{\mathbb{R}} \\ & \Downarrow & \\ \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbb{C}^{\text{op}}} & \rightleftarrows & (\bar{\mathbb{R}}^{\mathbb{D}})^{\text{op}} : \mathcal{M}_* \end{array}$$

Profoncteur et structures du noyau



$$\begin{aligned} \mathbb{C}^{\text{op}} \times \mathbb{D} &\rightarrow \bar{\mathbb{R}} \\ &\Downarrow \\ \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbb{C}^{\text{op}}} &\iff (\bar{\mathbb{R}}^{\mathbb{D}})^{\text{op}} : \mathcal{M}_* \end{aligned}$$

Projective metric geometry of tropical nuclei: gap matrices, event loci, and order chambers

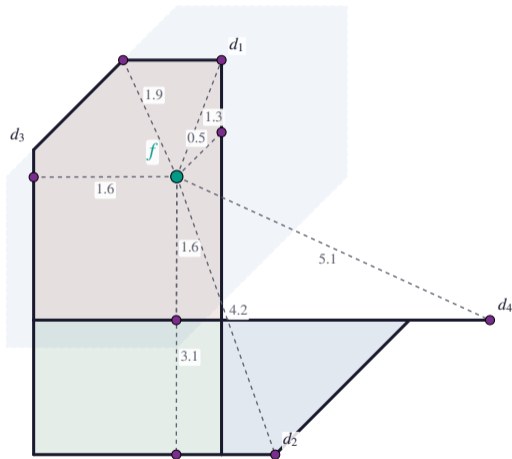
Juan Luis Gastaldi^{1†}, Samantha Jarvis^{2†}, Thomas Seiller^{3†},
John Terilla^{2,4†}

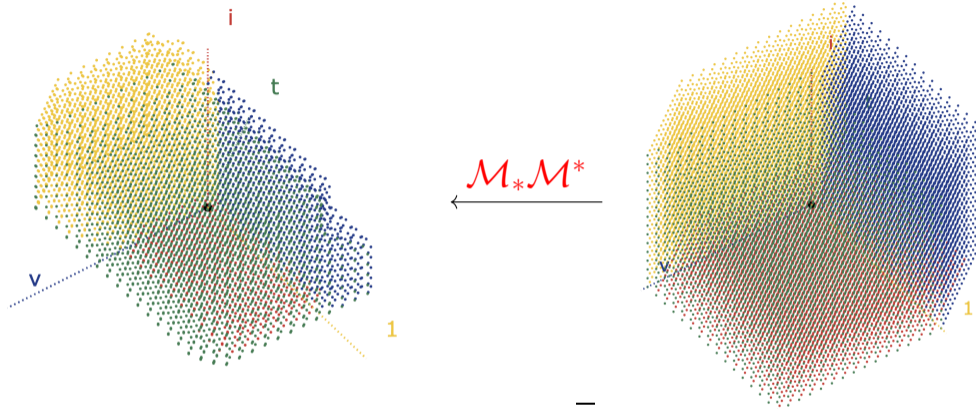
¹ETH Zurich, Zurich, Switzerland.

²Queens College, City University of New York, New York, NY, USA.

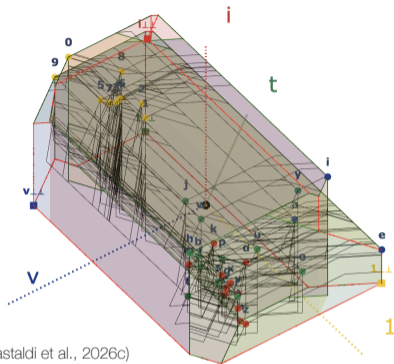
³CNRS, Paris, France.

⁴The Graduate Center, City University of New York, New York, NY, USA.





$$\begin{aligned}
 \mathbb{C}^{\text{op}} \times \mathbb{D} &\rightarrow \bar{\mathbb{R}} \\
 &\Downarrow \\
 \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbb{C}^{\text{op}}} &\iff (\bar{\mathbb{R}}^{\mathbb{D}})^{\text{op}} : \mathcal{M}_*
 \end{aligned}$$

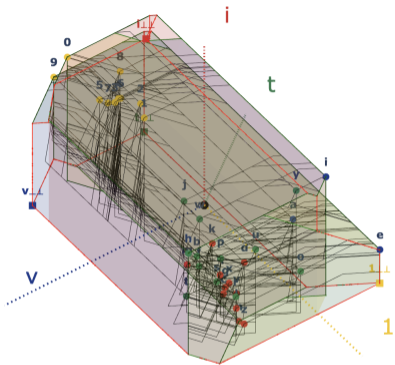


(Gastaldi et al., 2026c)

$$\begin{array}{ccc}
 \mathbf{C}^{\text{op}} \times \mathbf{D} & \rightarrow & \bar{\mathbb{R}} \\
 & \Downarrow & \\
 \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbf{C}^{\text{op}}} & \iff & (\bar{\mathbb{R}}^{\mathbf{D}})^{\text{op}} : \mathcal{M}_*
 \end{array}$$

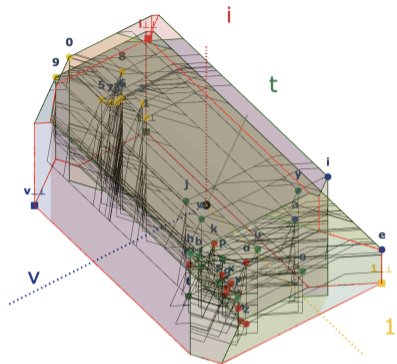
A red curved arrow points from the top-left node of the network diagram to the \mathcal{M}^* term in the equation above.

Géométrie du noyau



- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polhyedriques
- ◇ Invariant par jauge externe

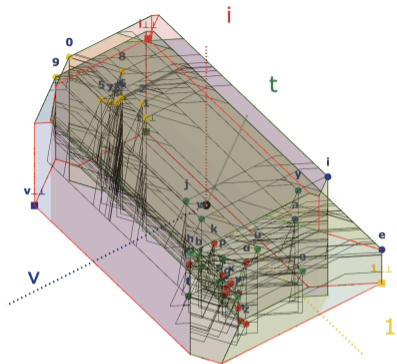
Géométrie du noyau



$$\mathcal{M} = \mathbf{C}^{\text{op}} \times \mathbf{D} \rightarrow \bar{\mathbb{R}}$$
$$\bar{\mathbb{R}} = ([-\text{inf}, \text{inf}], \leq, +)$$

- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polhyedriques
- ◇ Invariant par jauge externe

Géométrie du noyau



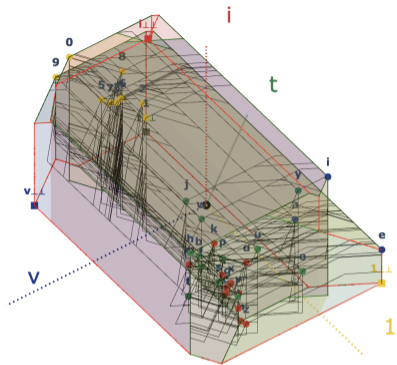
- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polhyedriques
- ◇ Invariant par jauge externe

$$\mathcal{M} = \mathbf{C}^{\text{op}} \times \mathbf{D} \rightarrow \bar{\mathbb{R}}$$
$$\bar{\mathbb{R}} = ([-\text{inf}, \text{inf}], \leq, +)$$

$$\mathcal{M}^* f(d) := \min_{c \in \mathbf{C}} (\mathcal{M}(c, d) - f(c))$$

$$\mathcal{M}^* g(c) := \min_{d \in \mathbf{D}} (\mathcal{M}(c, d) - g(d))$$

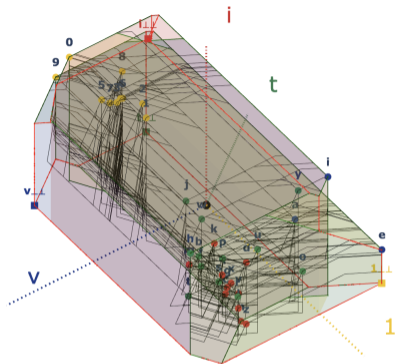
Géométrie du noyau



$$(f, g) \in \text{Nuc}(\mathcal{M}) \implies (f + \lambda, g - \lambda) \in \text{Nuc}(\mathcal{M})$$

- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polhyedriques
- ◇ Invariant par jauge externe

Géométrie du noyau

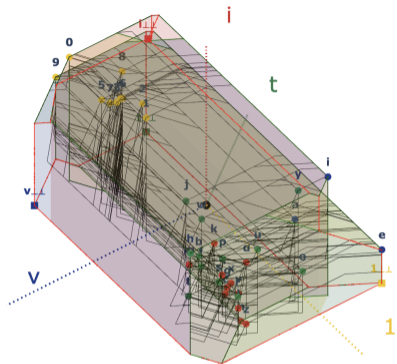


$$(f, g) \in \text{Nuc}(\mathcal{M}) \implies (f + \lambda, g - \lambda) \in \text{Nuc}(\mathcal{M})$$

La projectivisation $\mathbb{P}\text{Nuc}(\mathcal{M})$ est un
espace compact polyédrique

- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polyédriques
- ◇ Invariant par jauge externe

Géométrie du noyau

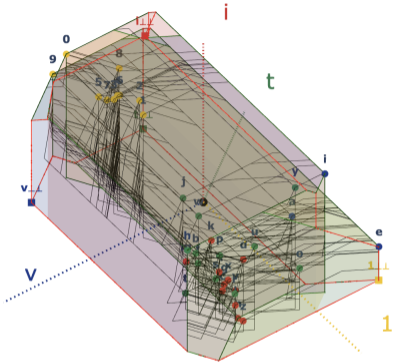


- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polhyedriques
- ◇ Invariant par jauge externe

$$d_C(f, f') = \max_c(f(c) - f'(c)) - \min_c(f(c) - f'(c))$$

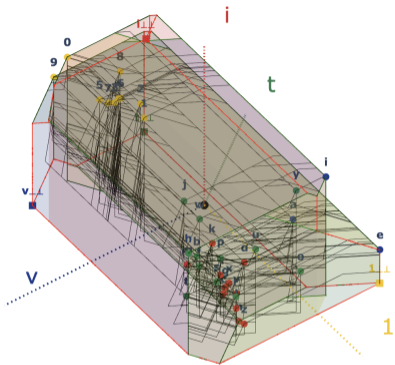
Géométrie du noyau

- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ **Décomposition en cellules polhyedriques**
- ◇ Invariant par jauge externe



Cellules définies par les couples (c, d) pour lesquelles $f(c) + g(d) = M(c, d)$

Géométrie du noyau



$$\delta^{f,g}(c, d) = \mathcal{M}(c, d) - f(c) - g(d)$$

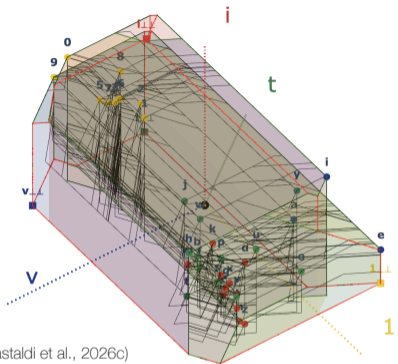
- ◇ Action de jauge
- ◇ Métrique projective tropicale interne
- ◇ Décomposition en cellules polhyedriques
- ◇ Invariant par jauge externe

A CALCULUS OF TYPES IN ISBELL NUCLEI

JUAN LUIS GASTALDI, SAMANTHA JARVIS, THOMAS SEILLER, AND JOHN TERILLA

ABSTRACT. We identify two constructions from different mathematical traditions. In linear logic and realisability, logical types are generated rather than fixed in advance: one begins with a universe of realisers equipped with execution, uses orthogonality to test their interactions, and takes types to be the biorthogonally closed subsets. In enriched Isbell duality, a quantitative relation induces an adjunction whose fixed points form a category, its nucleus. These constructions proceed by different means; we show that, in the present setting, they produce the same objects.

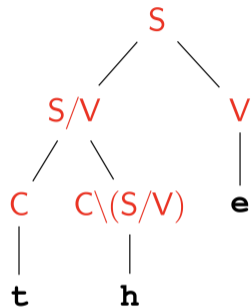
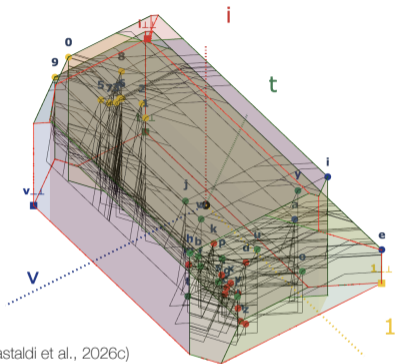
The shared datum is minimal: an associative product, called execution, and a real-valued measurement, with no compatibility assumed between them. The failure of the measurement to be additive is at once the relation defining orthogonality and the quantitative relation whose Isbell nucleus we form, and the types cut out by orthogonality are exactly the fixed points of the associated adjunction. The identification pays off in both directions. The most natural product of types fails to be associative; repairing this failure forces a different notion of type, sensitive to both sides of a composite, on which the induced product is associative and, when execution has units, carries two residuals. What emerges is a noncommutative Lambek calculus, derived directly from execution and orthogonality rather than imposed. In the reverse direction, each such type, read on the categorical side, generates a quantitative relation of its own, and with it a derived adjunction and a further generation of types; these derived types are again types of the original situation, computed by the residuals of the Lambek calculus. We also prove a coherence theorem for the threefold arrangements of this construction and, in the finite-dimensional case, give explicit formulas for the product.



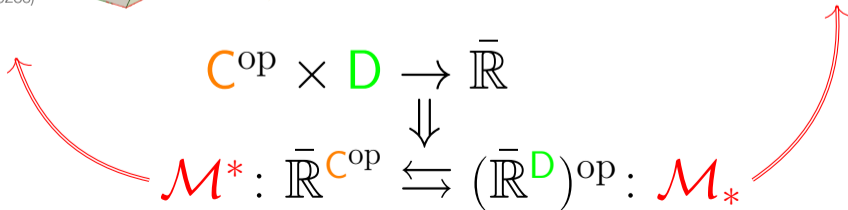
(Gastaldi et al., 2026c)

$$\begin{array}{ccc}
 \mathbf{C}^{\text{op}} \times \mathbf{D} & \rightarrow & \bar{\mathbb{R}} \\
 & \Downarrow & \\
 \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbf{C}^{\text{op}}} & \iff & (\bar{\mathbb{R}}^{\mathbf{D}})^{\text{op}} : \mathcal{M}_*
 \end{array}$$

Logique du noyau



(Gastaldi & Pellissier, 2021)
(Gastaldi et al., 2026a)



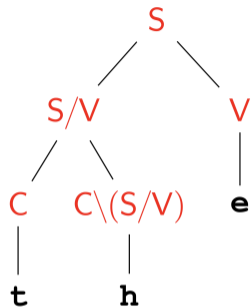
Définition: Polaire/Orthogonal - Girard, 2006

[É]étant donnée une fonction binaire $a, b \rightsquigarrow \langle a|b \rangle : A \times B \rightarrow C$ et un sous-ensemble $P \subset C$ (le « pôle »), on peut définir le *polaire* $X^\perp \subset B$ d'un sous-ensemble $X \subset A$ (resp. $Y^\perp \subset A$ d'un sous-ensemble $Y \subset B$) par :

$$X^\perp := \{y \in B : \forall x \in X, \langle a|b \rangle \in P\}$$

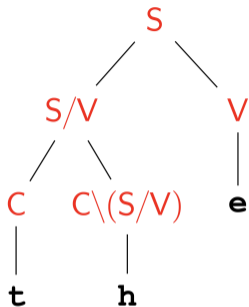
$$Y^\perp := \{x \in A : \forall y \in Y, \langle a|b \rangle \in P\}$$

- ◇ L'application « polaire » est décroissante:
 $X \subset X' \Rightarrow X'^\perp \subset X^\perp$.
- ◇ L'ensemble $\text{Pol}(A) \subset \mathcal{P}(A)$ des ensembles *polaires*, i.e., de la forme Y^\perp , est stable par intersections arbitraires. En particulier, A est polaire et $X^{\perp\perp}$ est le plus petit ensemble polaire contenant X .
- ◇ En conséquence, $X^{\perp\perp\perp} = X^\perp$.



Logique du Noyau

- ◇ Réalisabilité linéaire (Seiller, 2024)
- ◇ Produit tensoriel sur les types
- ◇ Implications à droite et à gauche

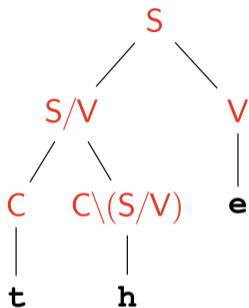


- ◇ Ensemble P
- ◇ $\bullet : P \times P \rightarrow P$ (opération associative)
- ◇ $[[\cdot, \cdot]]_m : P \times P \rightarrow \mathbf{R}$ ("measurement")

$$[[p_1 \bullet p_2, p_3]]_m + [[p_1, p_2]]_m = [[p_1, p_2 \bullet p_3]]_m + [[p_2, p_3]]_m$$

Logique du Noyau

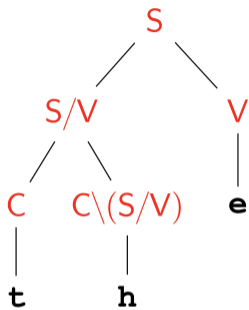
- ◇ Réalisabilité linéaire (Seiller, 2024)
- ◇ **Produit tensoriel sur les types**
- ◇ Implications à droite et à gauche



$$A \cdot B = \perp(\{a \cdot b \mid a \in A, b \in B\}^\perp)$$

Logique du Noyau

- ◇ Réalisabilité linéaire (Seiller, 2024)
- ◇ Produit tensoriel sur les types
- ◇ Implications à droite et à gauche

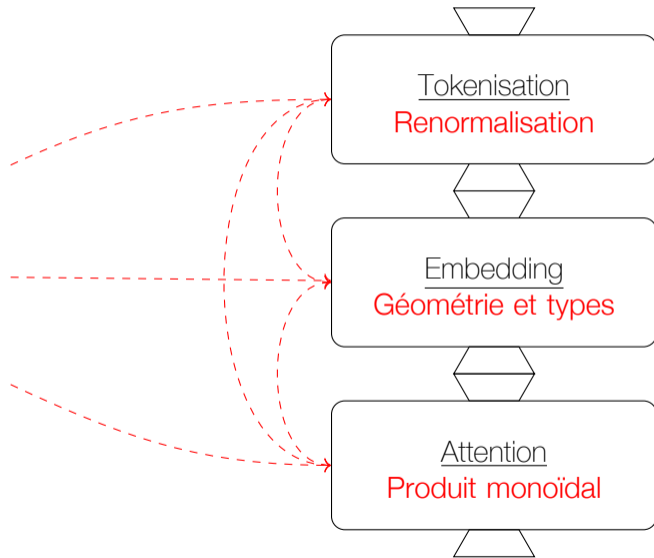
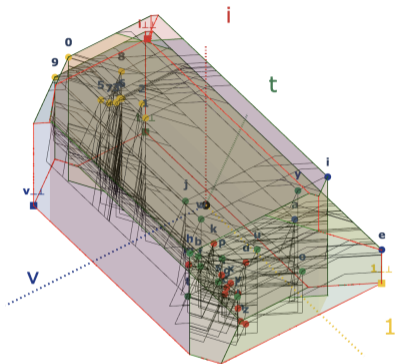


$$A \cdot B = \perp(\{a \cdot b \mid a \in A, b \in B\}^\perp)$$

$$A \multimap_r B = \{f \mid \forall a \in A, fa \in B\}$$

$$A \multimap_l B = \{f \mid \forall a \in A, af \in B\}$$

Explicabilité formelle



Intro: IA et sciences humaines

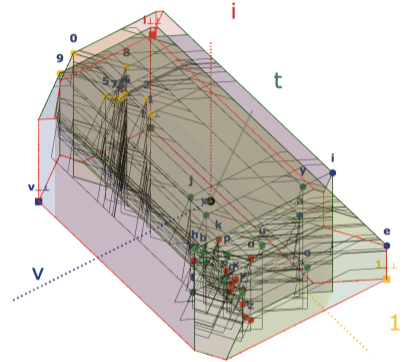
Explicabilité formelle

Interprétabilité théorique

Conclusion

Théorie

?

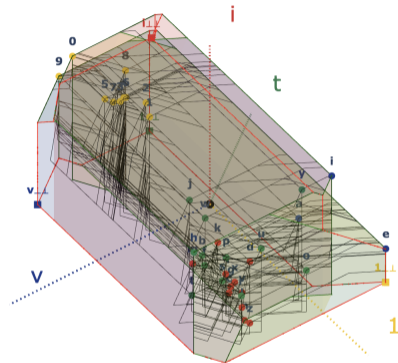


Hypothèse distributionnelle

Le contenu des unités linguistiques est déterminé par leur *distribution* dans un corpus.

$$C^{op} \times D \rightarrow \bar{\mathbb{R}}$$

Théorie



Hypothèse distributionnelle

Le contenu des unités linguistiques est déterminé par leur *distribution* dans un corpus.

$$C^{op} \times D \rightarrow \bar{R}$$



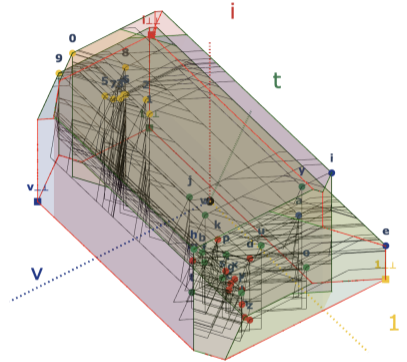
$$\bar{R}^{C^{op}} \Leftrightarrow (\bar{R}^D)^{op}$$

Théorie



Hypothèse structurale

Le contenu linguistique est l'effet d'une *structure* virtuelle dérivée des pratiques linguistiques dans une communauté.

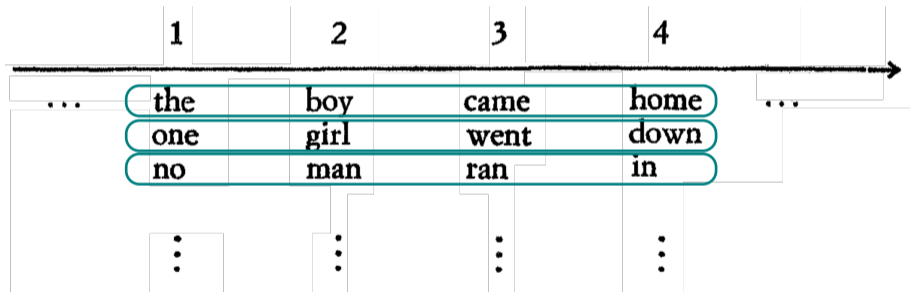


L'Hypothèse Structurale

“Il semble légitime en tous cas de poser a priori l'hypothèse qu'**à tout processus répond un système** qui permette de l'analyser et de le décrire au moyen d'**un nombre restreint de prémisses**. Il doit être possible de considérer tout processus comme composé d'**un nombre limité d'éléments** qui réapparaissent constamment dans de nouvelles combinaisons. On devrait pouvoir, en se fondant sur l'analyse du processus, **regrouper ces éléments en classes**, chaque classe étant définie par l'homogénéité de ses possibilités combinatoires, et pouvoir, à partir de ce classement préalable, **établir un calcul général exhaustif des combinaisons possibles**.”

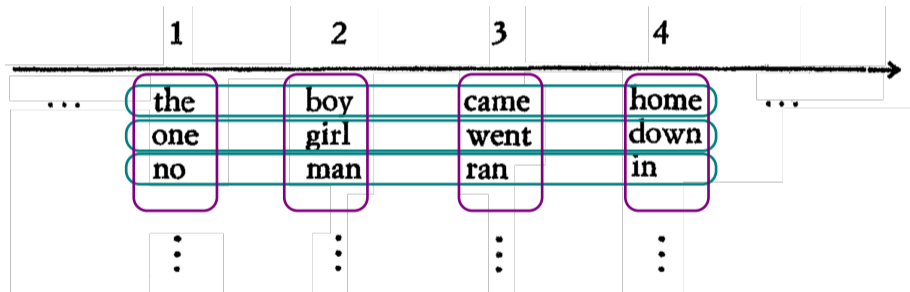
(Hjelmslev, *Prolégomènes à une théorie du langage*, p. 16)

Points fixes comme paradigmes



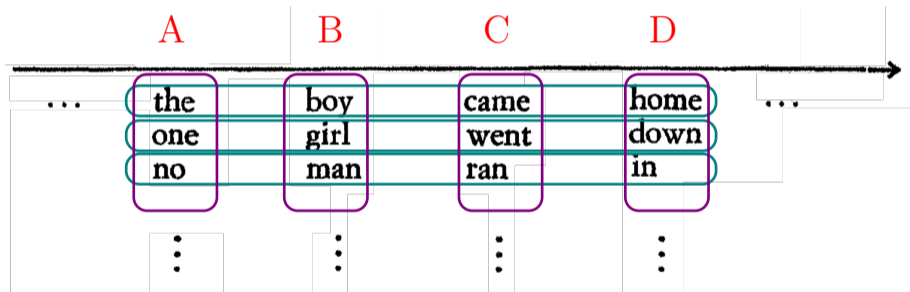
(Hjelmslev, 1971a)

Points fixes comme paradigmes



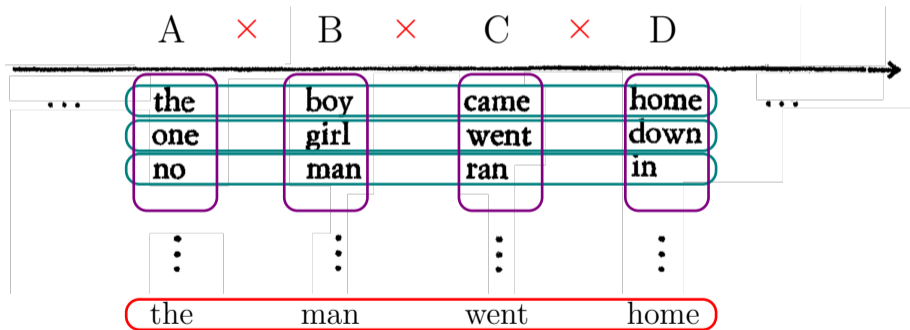
(Hjelmslev, 1971a)

Points fixes comme paradigmes



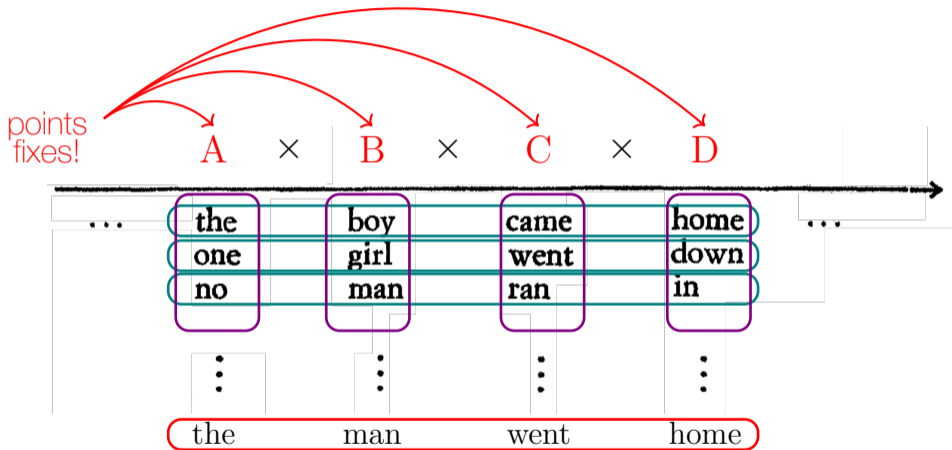
(Hjelmslev, 1971a)

Points fixes comme paradigmes



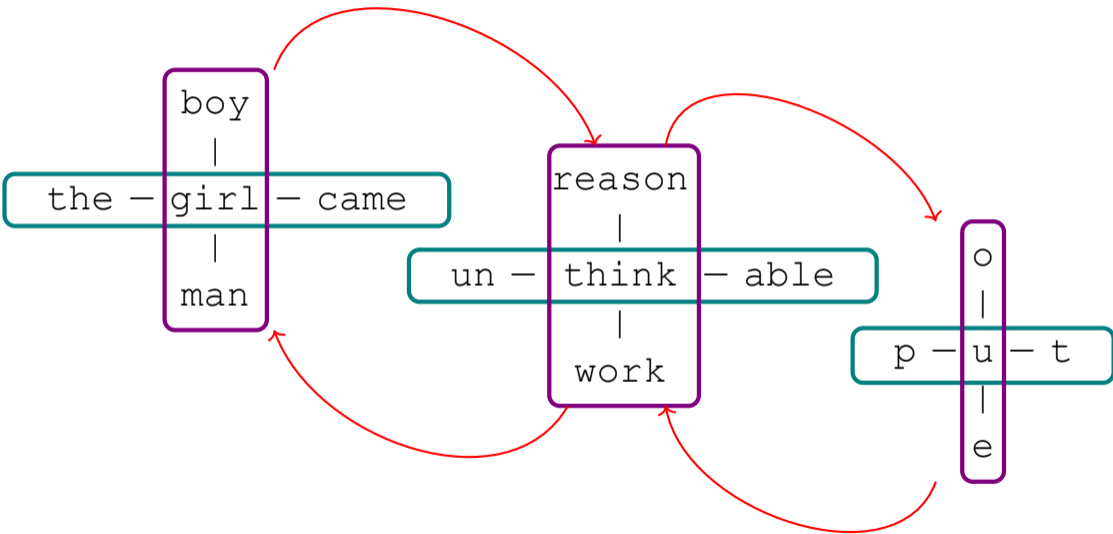
(Hjelmslev, 1971a)

Points fixes comme paradigmes



(Hjelmslev, 1971a)

Structure stratifiée



Points fixes comme traits distinctifs

	o	a	e	u	ə	i	l	ŋ	ʃ	ʃ̂	k	ʒ	ʒ̂	g	m	f	p	v	b	n	s	θ	t	z	ʒ	d	h	#	
1. Vocalic/Non-vocalic	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
2. Consonantal/Non-consonantal	-	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-
3. Compact/Diffuse	+	+	+	-	-	-	-	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
4. Grave/Acute	+	+	-	+	+	-	-	-	-	-	-	-	-	-	+	+	+	+	+	-	-	-	-	-	-	-	-	-	
5. Flat/Plain	+	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
6. Nasal/Oral								+	-	-	-	-	-	-	+	-	-	-	-	+	-	-	-	-	-	-	-	-	
7. Tense/Lax									+	+	+	-	-	-	-	+	+	-	-	-	+	+	+	-	-	-	-	+	-
8. Continuant/Interrupted									+	-	-	+	-	-	-	+	-	+	-	-	+	+	-	+	+	-	-	-	
9. Strident/Mellow											+	-	-	+	-	-	-	-	-	-	+	-	-	+	-	-	-	-	

(Jakobson et al., 1952)

Points fixes comme traits distinctifs

points
fixes!!!

points
fixes!!!

1. Vocalic/Non-vocalic
2. Consonantal/Non-consonantal
3. Compact/Diffuse
4. Grave/Acute
5. Flat/Plain
6. Nasal/Oral
7. Tense/Lax
8. Continuant/Interrupted
9. Strident/Mellow

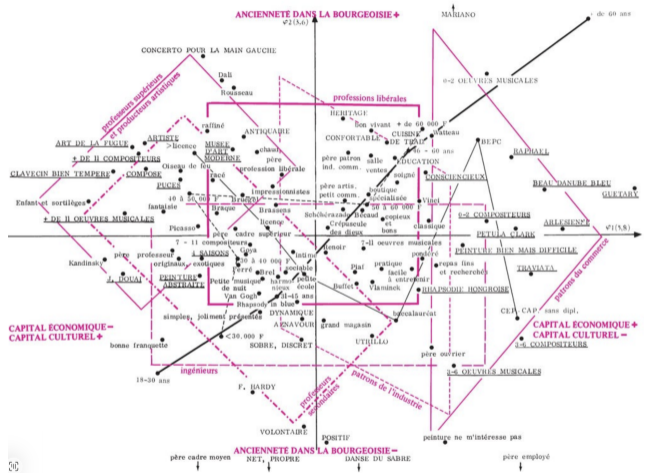
	o	a	e	u	ə	i	l	ŋ	ʃ	ʃ̂	k	z	ʒ	g	m	f	p	v	b	n	s	θ	t	z	ʒ	d	h	#	
1. Vocalic/Non-vocalic	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
2. Consonantal/Non-consonantal	-	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-
3. Compact/Diffuse	+	+	+	-	-	-	-	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
4. Grave/Acute	+	+	-	+	+	-	-	-	-	-	-	-	-	-	+	+	+	+	+	-	-	-	-	-	-	-	-	-	
5. Flat/Plain	+	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
6. Nasal/Oral								+	-	-	-	-	-	-	+	-	-	-	-	+	-	-	-	-	-	-	-	-	
7. Tense/Lax									+	+	+	-	-	-	-	+	+	-	-	-	+	+	+	-	-	-	-	+	-
8. Continuant/Interrupted									+	-	-	+	-	-	-	+	-	+	-	-	+	+	-	+	+	-	-	-	
9. Strident/Mellow											+	-	-	+	-	-	-	-	-	-	+	-	-	+	-	-	-	-	

(Jakobson et al., 1952)

Noyau et espace social

“S’il est vrai que [...] la **classe dominante** constitue un espace relativement autonome dont la **structure** est définie par la **distribution** entre ses membres des différentes espèces de capital, [...] on doit **retrouver ces structures** dans l’espace des styles de vie [...]. C’est ce que l’on a essayé d’établir en soumettant à **l’analyse des correspondances** l’ensemble des données recueillies.”

(Bourdieu, 1979)

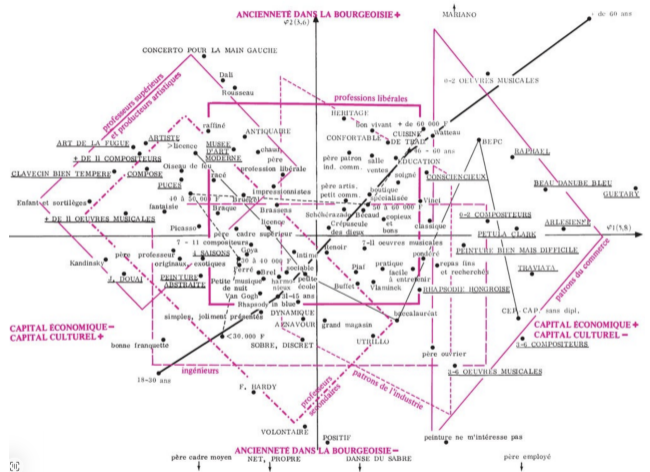


(Bourdieu, 1979)

Noyau et espace social

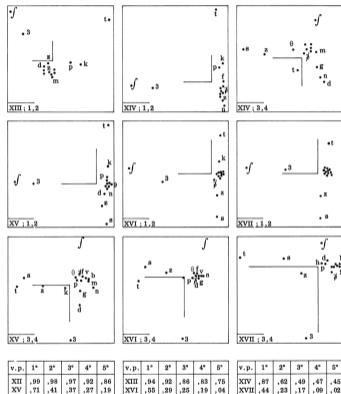
“Doing a **data analysis**, in **good mathematics**, is simply **searching eigenvectors**, all the science of it (the art) is just to find **the right matrix** to diagonalize”

(J.-P. Benzécri, 1973)



(Bourdieu, 1979)

Sciences humaines et sciences formelles



(J. P. Benzécri, 1976)

1) une voyelle neutre (amorphe), caractérisée par l'absence de chacune des propriétés $\beta, \phi, \chi, \lambda$: [a] ;

2) quatre types élémentaires de voyelles, chacun caractérisé par une seule propriété :

$$[e] = \phi, [a] = \beta, [v] = \chi, [z] = \lambda$$

3) six voyelles distinctes, chacune caractérisée par deux propriétés :

$$[o] = \phi\beta, [i] = \chi\phi, [u] = \chi\beta, [e] = \lambda\phi, [o] = \lambda\beta, [\partial] = \chi\lambda;$$

4) quatre voyelles combinées, chacune caractérisée par trois propriétés :

$$[e] = \chi\lambda\phi, [o] = \chi\lambda\beta, [u] = \phi\beta\chi, [\partial] = \phi\beta\lambda;$$

5) une voyelle polymorphe, caractérisée par les quatre propriétés considérées : la voyelle russe [sɨ] = $\phi\beta\chi\lambda$.

On obtient alors le diagramme de la figure 2.

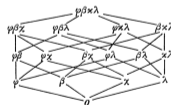
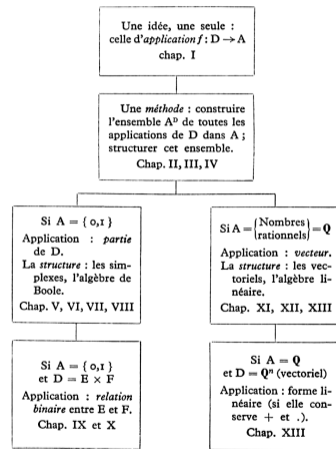


FIG. 2.

(Marcus, 1967)

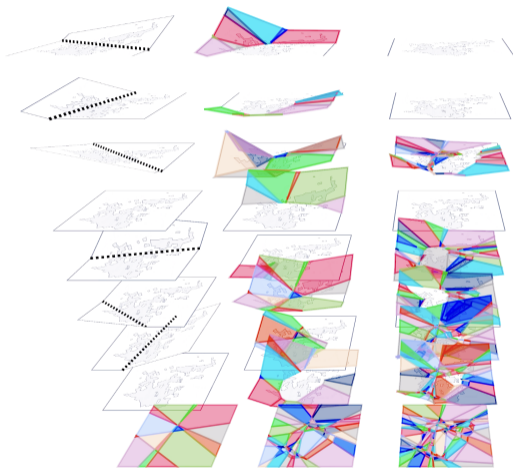
Organigramme



(Barbut, 1967)

Vers un structuralisme génératif

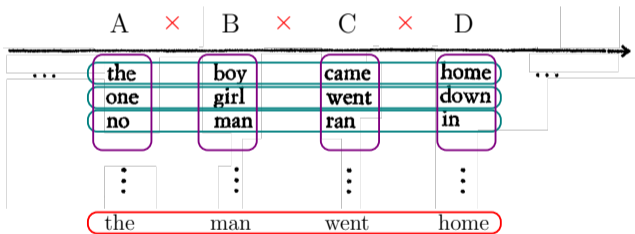
- ◇ Géométrie tropicale polyhédrique
- ◇ Réalisabilité et logique linéaire



(Welch Labs, 2025)

Vers un structuralisme génératif

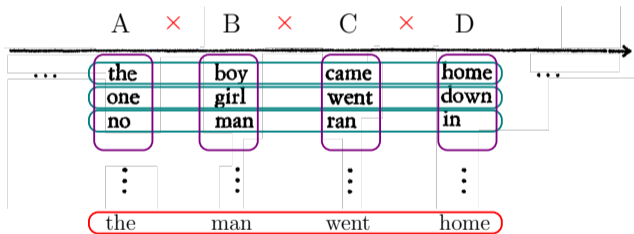
- ◇ Géométrie tropicale polyhédrique
- ◇ Réalisabilité et logique linéaire



(Hjelmslev, 1971a)

Vers un structuralisme génératif

- ◇ Géométrie tropicale polyédrique
- ◇ Réalisabilité et logique linéaire

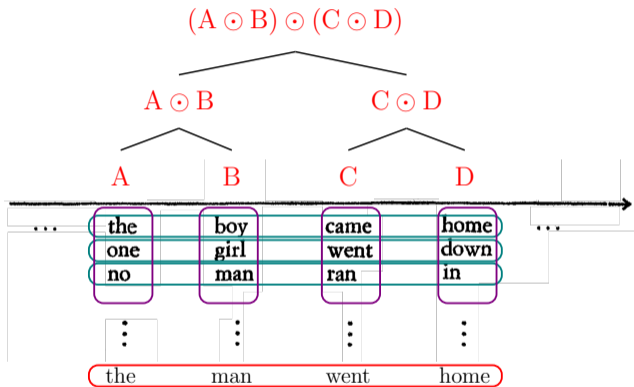


(Hjelmslev, 1971a)

$$(f \odot g)(x) := \sup_{cd=x} (f(c) + g(d) - M(c, d))$$

Vers un structuralisme génératif

- ◇ Géométrie tropicale polyhédrique
- ◇ Réalisabilité et logique linéaire



(Hjelmslev, 1971a)

$$(f \odot g)(x) := \sup_{cd=x} (f(c) + g(d) - M(c, d))$$

Intro: IA et sciences humaines

Explicabilité formelle

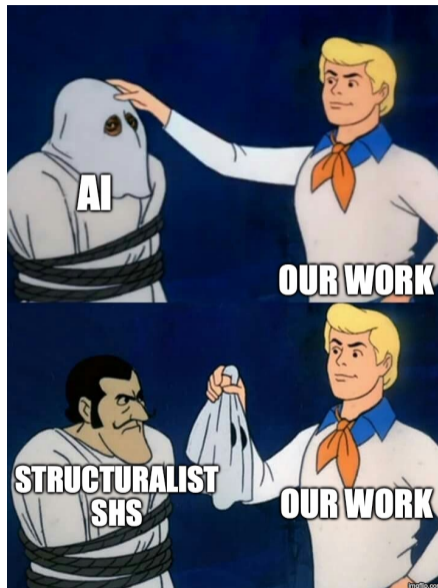
Interprétabilité théorique

Conclusion

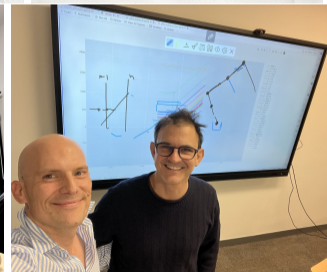
Conclusion



Conclusion



Collaborations



J. Terilla (CUNY), T.-D. Bradley (SandboxAQ), L. Pellissier (Paris-Est Créteil), Th. Seiller (CNRS), S. Jarvis (CUNY)

Reference Papers

- ◇ Gastaldi, J. L. (2021). Why Can Computers Understand Natural Language? *Philosophy & Technology*, 34(1), 149–214
- ◇ Gastaldi, J. L., & Pellissier, L. (2021). The Calculus of Language: Explicit Representation of Emergent Linguistic Structure through Type-Theoretical Paradigms. *Interdisciplinary Science Reviews*, 46(4), 569–590
- ◇ Bradley, T.-D., Gastaldi, J. L., & Terilla, J. (2024). The Structure of Meaning in Language: Parallel Narratives in Linear Algebra and Category Theory. *Notices of the American Mathematical Society*
- ◇ Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026c). Projective metric geometry of tropical nuclei: Gap matrices, event loci, and order chambers.
- ◇ Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026a). A calculus of types in isbell nuclei.

Références I

- Awodey, S. (2010). *Category theory* (2nd). Oxford University Press, Inc.
- Barbut, M. (1967). *Mathématiques et sciences humaines. tome i: Combinatoire et algèbre*. Presses Universitaires de France.
- Benzécri, J. P. (1976). Sur le codage réduit d'un vecteur de description en analyse des correspondances. *Les cahiers de l'analyse des données*, 1(2), 127–136.
- Benzécri, J.-P. (1973). *L'analyse des données. 2 L'analyse des correspondances* [Et coll.]. Bordas.
- Bourdieu, P. (1979). *La distinction: Critique sociale du jugement*. Éditions de Minuit.
- Bradley, T.-D., Gastaldi, J. L., & Terilla, J. (2024). The Structure of Meaning in Language: Parallel Narratives in Linear Algebra and Category Theory. *Notices of the American Mathematical Society*.
- Gastaldi, J. L. (2021). Why Can Computers Understand Natural Language? *Philosophy & Technology*, 34(1), 149–214.
- Gastaldi, J. L. (2024). Content from Expressions. The Place of Textuality in Deep Learning Approaches to Mathematics. *Synthese (under review)*.
- Gastaldi, J. L., & Pellissier, L. (2021). The Calculus of Language: Explicit Representation of Emergent Linguistic Structure through Type-Theoretical Paradigms. *Interdisciplinary Science Reviews*, 46(4), 569–590.
- Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026a). A calculus of types in isbell nuclei.
- Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026b). *Logical structures in \mathbb{R} -enriched adjunctions* [Working paper. Accessible at <https://www.giannigastaldi.com/assets/pdf/pubs/GastaldiJarvisEtAl2025.pdf>].
- Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026c). Projective metric geometry of tropical nuclei: Gap matrices, event loci, and order chambers.
- Geiger, A., Ibeling, D., Zur, A., Chaudhary, M., Chauhan, S., Huang, J., Arora, A., Wu, Z., Goodman, N., Potts, C., & Icard, T. (2025). Causal abstraction: A theoretical foundation for mechanistic interpretability.
- Girard, J.-Y. (2006). *Le point aveugle: Cours de logique. vers la perfection*. Editions Hermann.

Références II

- Goldberg, Y., & Levy, O. (2014). Word2vec explained: Deriving mikolov et al.'s negative-sampling word-embedding method. *CoRR*, *abs/1402.3722*.
- Hjelmslev, L. (1971a). La structure fondamentale du langage. In *Prolégomènes à une théorie du langage* [Prolégomènes à une théorie du langage] (pp. 177–231). Éditions de Minuit.
- Hjelmslev, L. (1971b). *Prolégomènes à une théorie du langage*. Éditions de Minuit.
- Jakobson, R., Fant, G. M., & Halle, M. (1952). *Preliminaries to speech analysis: The distinctive features and their correlates*. MIT Press.
- Levy, O., & Goldberg, Y. (2014). Neural word embedding as implicit matrix factorization. *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, 2177–2185.
- Marcus, S. (1967). *Introduction mathématique à la linguistique structurale*. Dunod.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *CoRR*, *abs/1310.4546*.
- Seiller, T. (2024). *Mathematical informatics* [Habilitation thesis].
- Sennrich, R., Haddow, B., & Birch, A. (2016). Neural machine translation of rare words with subword units. *Proceedings of the 54th Annual Meeting of the ACL*, 1715–1725.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc.

CAMS - EHESS

Séminaire Systèmes complexes en sciences sociales

2 Juin 2026 – Paris, France

Explicabilité formelle et interprétabilité théorique
des modèles neuronaux

Juan Luis (Gianni) Gastaldi

ETH zürich

www.giannigastaldi.com