

# A (de)tour through bias mitigation and analogy based ML

Défi HyAIAI

---

G. Alves (ex-PhD) V. Bhargava (ex-intern) **M. Couceiro** L. Galarraga  
E. Marquer (PhD) A. Napoli R. Querish (PhD) G. Zervakis (ex-PhD) ...

**Part I:** Bias mitigation...

**Part II:** Analogy based ML...

## **PART I: (Harmful) bias mitigation...**

# Harmful bias of ML Models

**ML models:** *designed* to have some bias that *guide* them in their tasks

## Expected bias:

Credit card default prediction	(good) <i>credit payment history</i>	↑
Hate speech prediction	(presence of) <i>offensive terms</i>	↑

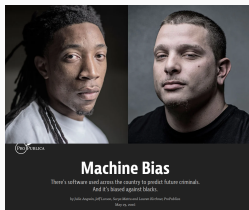
## Harmful bias:

Credit card default prediction	<i>ethnicity</i> (minority)	↓
Hate speech prediction	<i>language variant</i>	↓

Harmful bias lead to **unfair algorithmic decisions** & **discrimination**

**Discrimination:** “**unjust or prejudicial** treatment of different **categories of people**, especially, on the grounds of race, age, or sex”

# Motivation: unfair algorithmic decisions



COMPAS<sup>1</sup> (Tabular data)



Chatbot Tay<sup>2</sup> (Text)

Other Critical applications of algorithmic decisions: loan requests, job applications, Stop & Frisk, etc.

**Need of fairness:** Unfair outcomes not only affect human rights, but they undermine public trust in ML & AI.

---

<sup>1</sup> <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

<sup>2</sup> <https://www.bbc.com/news/technology-35902104>

# Addressing “fairness” in ML...

Based on **decision outcomes**, fairness can be assessed through:

- **Fairness metrics**: individual & group fairness, equal opportunity, demographic parity, equal accuracy, etc.
- **Process fairness**: model’s reliance on “sensitive features” (e.g., salient features such as race, age, or sex, . . . )

Two main approaches to tackle ML unfairness:

- **Enforce** fairness constraints while learning, e.g.:

$$P(y_{\text{pred}} \neq y_{\text{true}} | \text{race} = \text{Black}) = P(y_{\text{pred}} \neq y_{\text{true}} | \text{race} = \text{White})$$

**Drawback:** Complexity, “fairness overfitting”

- **Exclude** sensitive/salient features

**Drawback:** Decreased accuracy!

# Addressing “fairness” in ML...

Based on **decision outcomes**, fairness can be assessed through:

- **Fairness metrics**: individual & group fairness, equal opportunity, demographic parity, equal accuracy, etc.
- **Process fairness**: model’s reliance on “sensitive features” (e.g., salient features such as race, age, or sex, . . . )

Two main approaches to tackle ML unfairness:

- **Enforce** fairness constraints while learning, e.g.:

$$P(y_{\text{pred}} \neq y_{\text{true}} | \text{race} = \text{Black}) = P(y_{\text{pred}} \neq y_{\text{true}} | \text{race} = \text{White})$$

**Drawback:** Complexity, “fairness overfitting”

- **Exclude** sensitive/salient features

**Drawback:** Decreased accuracy!

**Fairness through unawareness...**



# FixOut (Fairness through eXplanations and feature dropOut)<sup>3</sup>

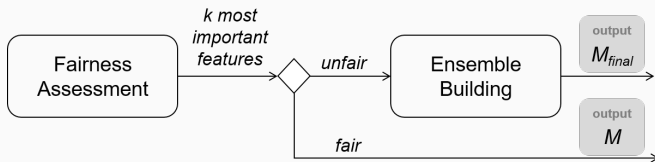
**Goal:** reduce model's dependence on sensitive/salient features **while** keeping (or improving) its performance

**Fair Model:** if its outcomes do not depend on sensitive features

## FixOut: Human-centered approach to deal process fairness

**Input:** model  $M$ , dataset  $D$ , sensitive features  $F$ , explanation method  $E$

**Output:**  $M$  if fair, **otherwise** a fairer and more accurate  $M_{final}$



<sup>3</sup><https://fixout.loria.fr/>

## Example: FixOut with LIME (RF on German)<sup>4</sup>

Original		Ensemble	
Feature	Contrib.	Feature	Contrib.
<b>foreignworker</b>	2.664899	otherinstallmentplans	-1.487604
otherinstallmentplans	-1.354191	housing	-1.089726
housing	-1.144371	savings	0.679195
savings	0.984104	duration	-0.483643
property	-0.648104	<b>foreignworker</b>	0.448643
purpose	-0.415498	property	-0.386355
existingchecking	0.371415	credithistory	0.258375
<b>telephone</b>	0.311451	job	-0.252046
credithistory	0.263366	existingchecking	-0.21358
duration	-0.223288	residencesince	-0.138818

**Result:**  $M_{\text{final}}$  is “fairer” & at least as accurate (from 0.783 to 0.786)

<sup>4</sup> Bhargava, et al. LimeOut: An Ensemble Approach to Improve Process Fairness. PKDD/ECML Workshop XKDD 2020: 475-491

**Q: Impact of FixOut on w.r.t. fairness metrics?**

# What about Fairness metrics?

**Idea:** Separate instances into two groups w.r.t. a sensitive feature  
**E.g.:** **Non-white people** (unprivileged) **versus white people** (privileged)

**Demographic Parity (DP)**<sup>5</sup> :

$$DP = P(\hat{y} = pos | D = unp) - P(\hat{y} = pos | D = priv)$$

**Equal Opportunity (EO)**<sup>6</sup>:  $EO = \frac{TP_{unp}}{TP_{unp} + FN_{unp}} - \frac{TP_{priv}}{TP_{priv} + FN_{priv}}$

**Predictive Equality (PE)**<sup>7</sup>:  $PE = \frac{FP_{unp}}{FP_{unp} + TP_{unp}} - \frac{FP_{priv}}{FP_{priv} + TP_{priv}}$

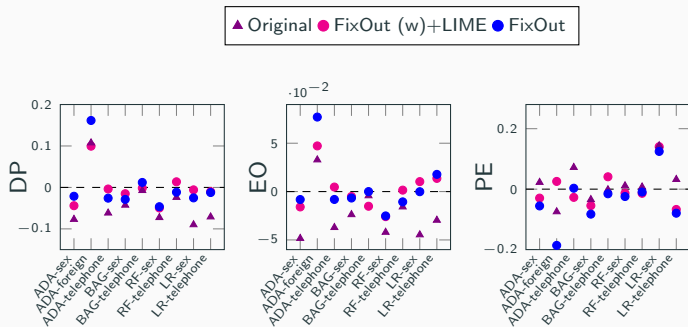
---

<sup>5</sup> Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 2017, 153–163.

<sup>6</sup> Zafar, et al. Fairness beyond disparate treatment & impact: Learning classification without disparate mistreat. WWW 2017.

<sup>7</sup> Alves, et al. Making ML models fairer through explanations: the case of LimeOut. AIST 2020.

# Fairness metrics



## German dataset: Privileged groups

- “status sex”: “male single”
- “telephone”: “yes” (registered under the customers name)
- “foreign worker”: “no”

**FixOut:** A human-centered approach to mitigate harmful bias

① On tabular data:

Use of explanations and **Control** through aggregation

**Automated** the choice of the most important features to be considered in the fairness assessment

② On textual data:

**How to** adapt *feature dropout* to *bag of words*.

**Reduced** unintended bias of ML models on textual data<sup>8</sup>

③ Adaptation to neural classifiers:

**Feature dropout** on the representations (embeddings)

---

<sup>8</sup> [Alves, et al.](#) Reducing Unintended Bias of ML Models on Tabular and Textual Data. DSAA 2021: 1-10.

**FixOut:** currently in **startup pre-maturation** (SATT & Incubateur Lorrain) to be followed by INRIA Startup Studio

**Further results:** Statistical approach (Hilbert-Schmidt IC) to detecting sensitive attributes (data-driven)<sup>9</sup>

**Refine-LM:** Reinforcement Learning for Harmful Bias Mitigation in LL Models

**Further results:** Portable bias filtering mechanism that is

- easy to train,
- adjustable to a multiple language models,
- adaptable to different bias contexts (gender, ethnicity, religion, etc.)

---

<sup>9</sup> [Pelegina, et al.](#) A statistical approach to detect sensitive features in a group fairness setting. CoRR abs/2305.06994 (2023).

**FixOut:** currently in **startup pre-maturation** (SATT & Incubateur Lorrain) to be followed by INRIA Startup Studio

**Further results:** Statistical approach (Hilbert-Schmidt IC) to detecting sensitive attributes (data-driven)<sup>9</sup>

**Refine-LM:** Reinforcement Learning for Harmful Bias Mitigation in LL Models

**Further results:** Portable bias filtering mechanism that is

- easy to train,
- adjustable to a multiple language models,
- adaptable to different bias contexts (gender, ethnicity, religion, etc.)

---

<sup>9</sup> [Pelegina, et al.](#) A statistical approach to detect sensitive features in a group fairness setting. CoRR abs/2305.06994 (2023).



## Part II: Analogy based ML...

## Example: Analogical proportion (a is to b as c is to d)



**Analogies** simultaneously exploit **similarities** and **dissimilarities**

**3 key cognitive processes:** **Abstraction**, **Inference** and **Creativity**

**Detecting/mining analogies:** Given  $a$ ,  $b$ ,  $c$ , and  $d$ ,

- is  $(a, b, c, d)$  a valid analogy?

**Solving analogies:** Given  $a$ ,  $b$ ,  $c$ ,

- find  $x$  s.t.  $(a, b, c, x)$  a valid analogy

**Reasoning and integrating analogical reasoning (AR):**

- Depending on the concrete application and ML&AI task

## Example: Analogical proportion (a is to b as c is to d)



**Analogies** simultaneously exploit **similarities** and **dissimilarities**

**3 key cognitive processes:** **Abstraction**, **Inference** and **Creativity**

**Detecting/mining analogies:** Given  $a$ ,  $b$ ,  $c$ , and  $d$ ,

- is  $(a, b, c, d)$  a valid analogy?

**Solving analogies:** Given  $a$ ,  $b$ ,  $c$ ,

- find  $x$  **s.t.**  $(a, b, c, x)$  a valid analogy

**Reasoning and integrating analogical reasoning (AR):**

- Depending on the concrete application and ML&AI task

# Different views on analogies

**Axiomatic:** As 4-ary relations satisfying certain postulates

**Examples:** reflexivity, (certain) permutations, etc.

**Relational:**  $R(a, b, c, d) \equiv P(P_1(a, b), P_1(c, d))$ , for  $P, P_1$  predicates

**Example:**  $R(\text{wine}, \text{France}, \text{beer}, \text{Germany})$

**Functional:**  $R(a, b, c, d)$  if  $b = T(a)$  and  $d = T(c)$ , for some  $T$

**Example:**  $R(\text{go}, \text{went}, \text{make}, \text{made})$

**Model Theoretic:** Relying on structural transformations and “rewriting”

**Examples:** *Structure mapping theory* and *Justifications*

**NB:** Different ways to define analogies depending on the data, the underlying structure and the task at hand...

# Different views on analogies

**Axiomatic:** As 4-ary relations satisfying certain postulates

**Examples:** reflexivity, (certain) permutations, etc.

**Relational:**  $R(a, b, c, d) \equiv P(P_1(a, b), P_1(c, d))$ , for  $P, P_1$  predicates

**Example:**  $R(\text{wine}, \text{France}, \text{beer}, \text{Germany})$

**Functional:**  $R(a, b, c, d)$  if  $b = T(a)$  and  $d = T(c)$ , for some  $T$

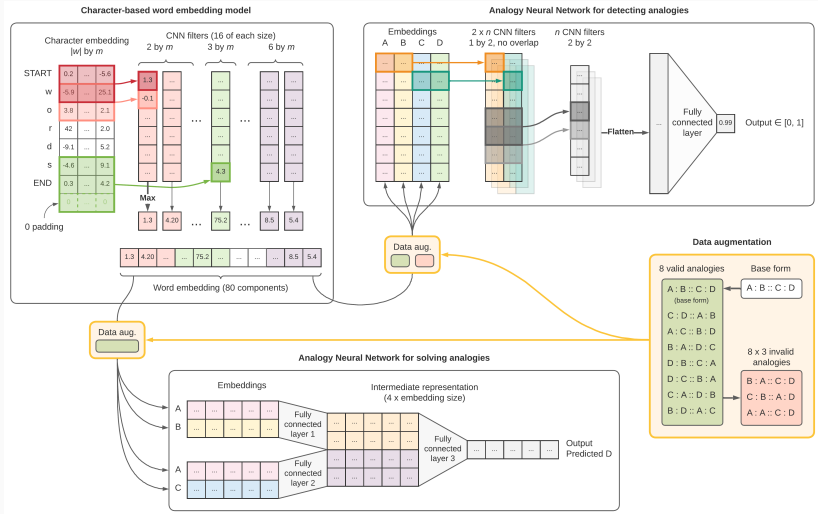
**Example:**  $R(\text{go}, \text{went}, \text{make}, \text{made})$

**Model Theoretic:** Relying on structural transformations and “rewriting”

**Examples:** *Structure mapping theory* and *Justifications*

**NB:** Different ways to define analogies depending on the data, the underlying structure and **the task** at hand...

# Example: detecting and solving morphological analogies



Recently: solving morphological analogies through generation [C22b]

ANNA: <https://anna.loria.fr/>

## Some application domains

- NLP & translation [L03, S05,M20,A21a,M22,C22]
- Classification & recom. [B07,B17, C17,C18,C20b,C23b]
- CB & Machine reasoning [F89,G83,L21,L19a,L19b,L21,M21]
- Transfer learning [B19,C20a, A21b,F23,M23]
- VisualQA, ScholasticAP, TSV, Explainability [S15,P19,Z22,H20]
- ...and even humor: pun and meme generation (WIP)



# Target Sense Verification (TSV)<sup>10</sup>

**Here:** Target Sense Verification (TSV)

**PhD:** Georgios Zervakis (defended March, 2023)

*Enriching large language models with semantic lexicons and analogies*

Intended Sense	Target Sense	
Context	Definition	Hypernyms
deliver the package to my <u>home</u>	where you live at a particular time	residence, abode
<u>home</u> is where the heart is	where you live at a particular time	residence, abode

**Question:** Is the *intended sense* of **home** the same as in the *target sense*?

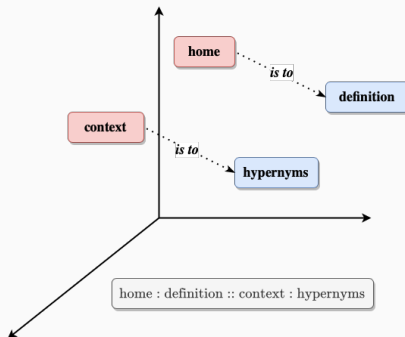
<sup>10</sup>Breit et al. (2021) WiC-TSV: An evaluation benchmark for target sense verification of words in context. EACL



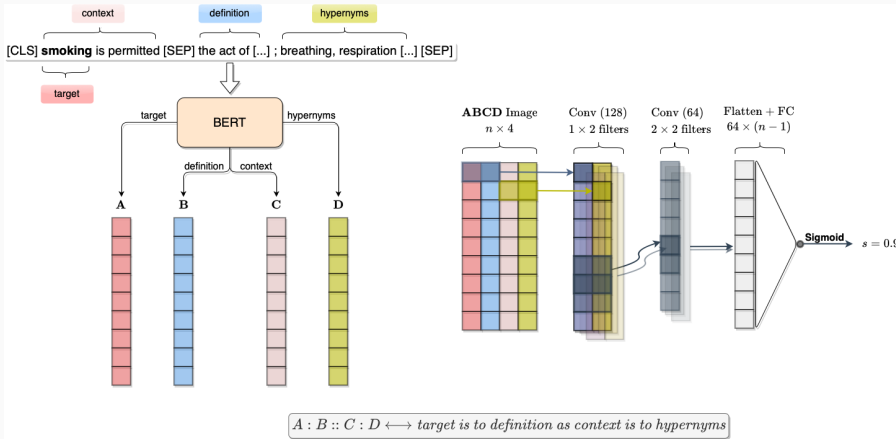
# TSV as analogy detection

Intended Sense		Target Sense
<b>Context</b> deliver the package to my <u>home</u>	True	<b>Definition</b> where you live at a particular time
<b>Context</b> <u>home</u> is where the heart is	False	<b>Definition</b> where you live at a particular time
		<b>Hypernyms</b> residence, abode

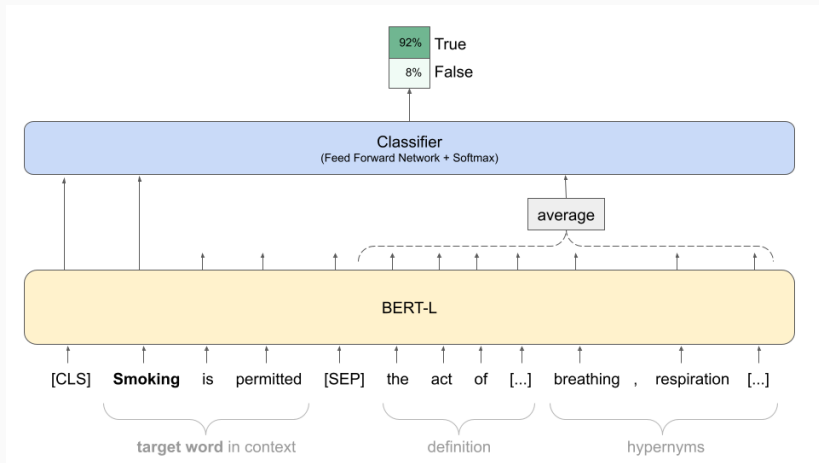
**Idea:** Formulate the question as an analogy and check whether it is valid.



# AB4TSV architecture



# Baselines (Breit et al.: WiC-TSV 2021)

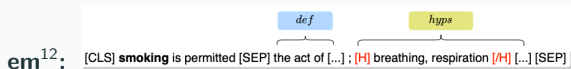
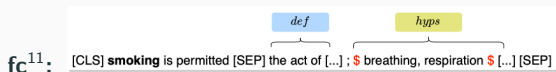
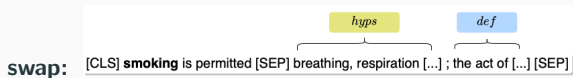


BERT + [CLS] + target word + average(definition, hypernyms)  $\Rightarrow$  Classifier

# Choice of input encoding and analogy relation



## Alternative input encoding operations

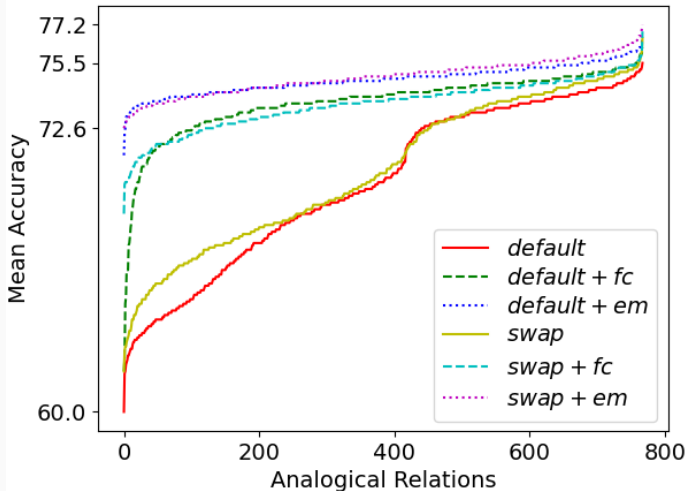


<sup>11</sup>Huang et al. (2019) GlossBERT: BERT for word sense disambiguation with gloss knowledge. EMNLP-IJCNLP.

<sup>12</sup>Baldini Soares et al. (2019) Matching the blanks: Distributional similarity for relation learning. ACL.

# Impact of input encoding

Mean accuracy: 6 input encodings  $\times$  768 analog. relations  $\times$  4 random seeds



# Analogical relation optimization results

Encoding	Analogy	Dev Acc	Dev F1
default	<i>cls : descr :: cls : ctx</i>	74.5 $\pm$ 0.015	77.0 $\pm$ 0.016
default+fc	<i>cls : def :: ctx : cls</i>	74.9 $\pm$ 0.010	77.3 $\pm$ 0.006
default+em	<i>tgt : descr :: cls : def</i>	75.4 $\pm$ 0.027	<b>77.8</b> $\pm$ 0.023
swap	<i>def : cls :: cls : ctx</i>	75.4 $\pm$ 0.016	77.7 $\pm$ 0.016
swap+fc	<i>def : ctx :: cls : hyps</i>	<b>75.8</b> $\pm$ 0.013	77.7 $\pm$ 0.013
swap+em	<i>hyps : def :: cls : ctx</i>	<b>75.8</b> $\pm$ 0.017	77.7 $\pm$ 0.012
Baselines			
default		74.0 $\pm$ 0.014	76.9 $\pm$ 0.007
default+fc		73.9 $\pm$ 0.018	76.3 $\pm$ 0.018
default+em	HyperBert3	73.1 $\pm$ 0.031	75.2 $\pm$ 0.032
swap		73.8 $\pm$ 0.015	76.3 $\pm$ 0.015
swap+fc		73.5 $\pm$ 0.011	75.6 $\pm$ 0.013
swap+em		74.4 $\pm$ 0.011	75.7 $\pm$ 0.024

**NB1:** AB4TSV >> Baselines

**NB2:** [CLS] token matters

## Comparative results

Approach	Test Acc	Test F1
CTLR <sup>13</sup>	78.3	78.5
V <sup>14</sup>	71.9	76.2
BERT <sub>Base</sub> <sup>15</sup>	76.6	78.2
BERT <sub>Large</sub> <sup>15</sup>	76.3	77.8
FastText <sup>15</sup>	53.4	63.4
AB4TSV+swap+em	75.7	77.5
AB4TSV+swap+fc	<b>78.6</b>	<b>79.8</b>
AB4TSV <sub>pi</sub> +default+em	<b>78.6</b>	79.4
U-dBERT <sup>15</sup>	61.2	51.3
U-BERT <sup>15</sup>	60.5	51.9
MIRRORWIC <sup>16</sup>	73.7	—

<sup>13</sup>Moreno et al. (2021) CTLR@WiC-TSV

<sup>14</sup>Vandenbussche et al. (2021) SemDeep-6

<sup>15</sup>Breit et al. (2021) WiC-TSV

<sup>16</sup>Liu et al. (2021) MirrorWiC

## Analogy and BERT for target sense verification (AB4TSV)

- ① Combining LLMs with analogy classifiers improves results on TSV
- ② Marking the input text with special characters can boost the performance.
- ③ Integrating the properties of analogies offers gains in interpretability.
- ④ AB4TSV shows OOD generalization and transfer learning capabilities.





<https://at2ta.loria.fr/>

PRCE Axis E.2, CES 23: Intelligence artificielle et science des données

## Partners:



# General objective and Challenges

**AT2TA General Objective:** propose an ML framework that integrates analogical reasoning (AR), easily adaptable to different real use cases.

**(C1) Bridging the gap between ML and KRR**

**(C2) Analogy modeling and representation learning for AR**

**(C3) AR adaptation across domains**

**(C4) Platform for multimodal/multi-domain AR**

## Events & dissemination:

- Annual workshops with proceedings:  
IARML@IJCAI-ECAI 2022 & ATA@ICCBR 2022 (both published)  
IARML@IJCAI 2023 & ATA@ICCBR 2023 (upcoming!)
- [Springer special issues](#) (yearly):  
*Analogies: from Mathematical Foundations to Applications and Interactions with ML and AI* (S722: Analogical reasoning)  
in [Annals of Mathematics and Artificial Intelligence \(AMAI\)](#)
- [Shared Tasks](#) (upcoming)
- Other actions (to discuss)

*Merci de votre attention!*

*Thank you for your attention!*

# Selected References

- [A22] Ch. Antić. Analogical proportions. *AMAI* 90:6 (2022) 595–644
- [A21a] S. Alsaidi, *et al.* A neural approach for detecting morphological analogies. *DSAA21*, 1-10
- [A21b] S. Alsaidi, *et al.* On the transferability of neural models of morphological analogies. *AIMLAI@ECML-PKDD21*, 76–89
- [B23] F.Badra, *et al.* Some Perspectives on Similarity Learning for Case-Based Reasoning and Analogical Transfer. To appear.
- [B19]b Z. Bouraoui, *et al.* From Shallow to Deep Interactions Between Knowledge Representation, Reasoning and Machine Learning (Kay R. Amel group), 2019. <http://arxiv.org/abs/1912.06612>
- [B17] M. Bounhas, *et al.* Analogy-based classifiers for nominal or numerical data. *IJAR* 91 (2017) 36–55
- [B07] S. Bayouhd, *et al.* Learning by analogy: A classification rule for binary and nominal data. *IJCAI07*, 678–683
- [C23] M. Couceiro, E. Lehtonen. Galois theory for analogical classifiers. In press *AMAI*, 2023. 21 p.
- [C22] K. Chan, *et al.* Solving Morphological Analogies Through Generation. *IARML@IJCAI-ECAI 2022*, 29–39
- [C20a] A. Cornuéjols, *et al.* Transfer learning by learning projections from target to source. *IDA20*, 119–131
- [C20b] M. Couceiro, *et al.* When nominal analogical proportions do not fail. *SUM20*, 68–83
- [C18] M. Couceiro, *et al.* Behavior of analogical inference w.r.t. Boolean functions. *IJCAI18*, 2057–2063
- [C17] M. Couceiro, *et al.* Analogy-preserving functions: A way to extend Boolean samples. *IJCAI17*, 1575–1581
- [G83] D. Gentner. Structure-Mapping: A Theoretical Framework for Analogy. *Cognitive Science* 7(2) (1983) 155–170
- [H16] N. Hug, *et al.* Analogical classifiers: A theoretical perspective. *ECAI16*, 689–697
- [H20] E. Hüllermeier. Towards Analogy-Based Explanations in ML. *CoRR* abs/2005.12800, 2020

## Selected References (cont.)

- [L21] J. Lieber, *et al.* When revision-based case adaptation meets analogical extrapolation. *ICCB21*, 156–170
- [L19a] J. Lieber, *et al.* Improving analogical extrapolation using case pair competence. *ICCB19*, 251–265
- [L19b] S. Lim, *et al.* Solving word analogies: A machine learning perspective. *ECSQARU19*, 238–250
- [L03] Y. Lepage. De l'analogie rendant compte de la commutation en linguistique. HdR Université Joseph Fourier - Grenoble I, 2003
- [M23] E. Marquer, *et al.* Less is Better: An Energy-Based Approach to Case Base Competence. To appear.
- [M22] E. Marquer, *et al.* A Deep Learning approach to solving morphological analogies. *ICCB22*, 159–174
- [M21] M. Mitchell. Abstraction and analogy-making in artificial intelligence, 2021
- [M20] P.-A. Murena, *et al.*. Solving analogies on words based on minimal complexity transformation. *IJCAI20*, 1848–1854
- [P19] J. Peyre, *et al.* Detecting unseen visual relations using analogies. *ICCV19*
- [S05] N. Stroppa, *et al.* Analogical learning and formal proportions: Definitions and methodological issues. Technical Report D004, ENST-Paris, 2005
- [Z22] G. Zervakis, *et al.* An Analogy based Approach for Solving Target Sense Verification. NLP1R 2022