

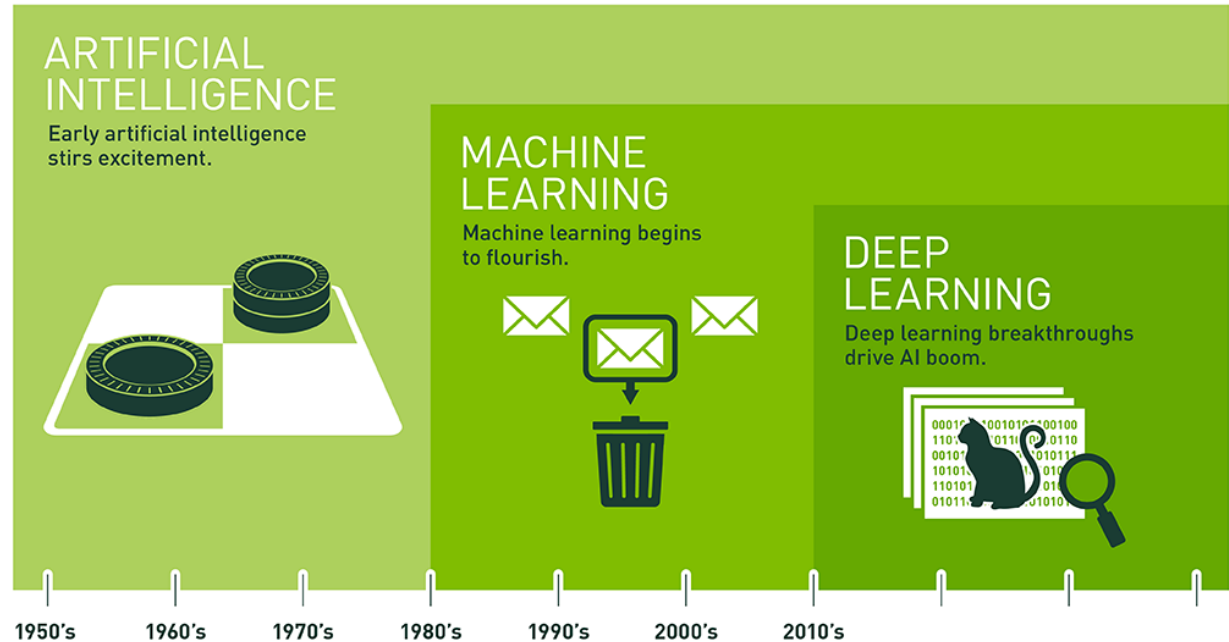
Explainable interpretations for the Entity Resolution task

Donatella Firmani

[donatella.firmani@u
niroma3.it](mailto:donatella.firmani@uniroma3.it)

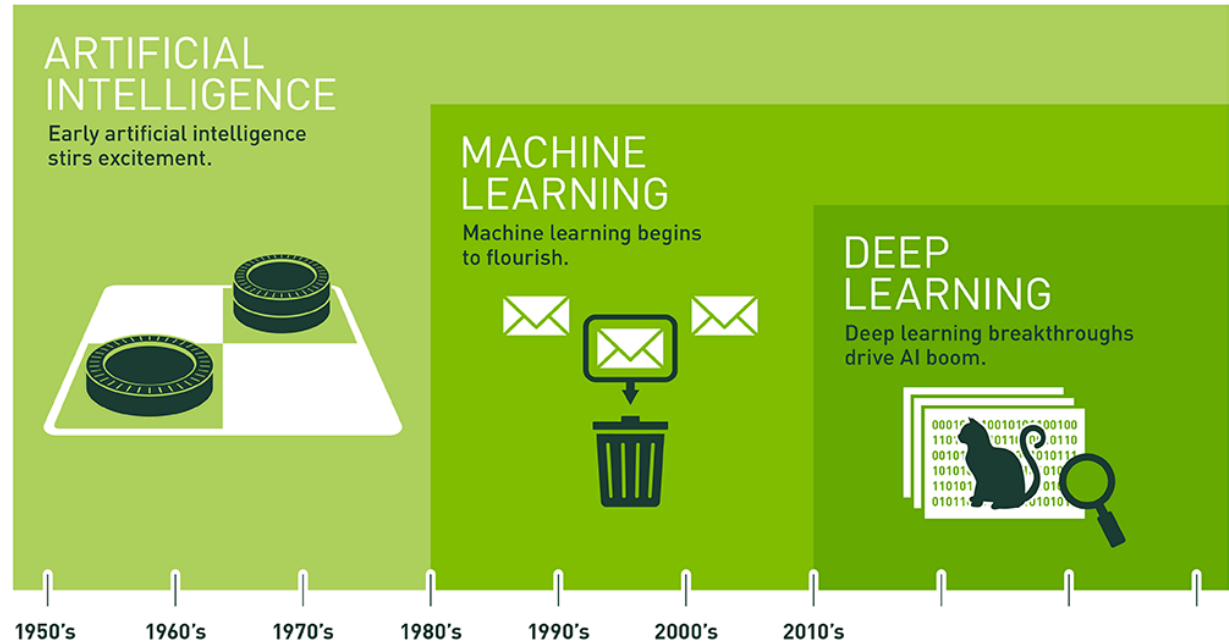
Big Data Seminars
2020

Brief History of AI



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

Brief History of AI



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger subsets of artificial intelligence.

Big Data



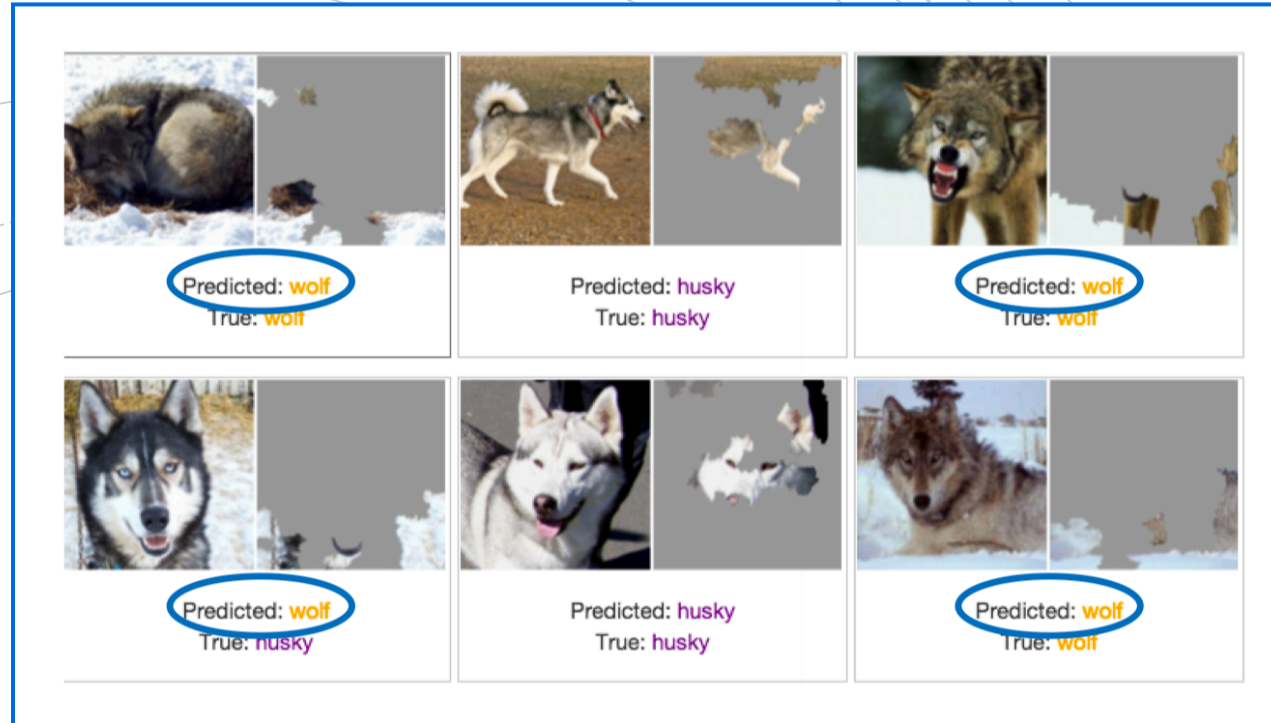
Data

- By relying on patterns in training data, machine learning can solve a specific task without using explicit instructions.
- Unprecedented accuracy in many application scenarios

Husky vs Wolf



- Only 1 prediction mistake: great accuracy
- Interested in “Why” questions, rather than “How Accurate”?



Or snow vs land?

Pixels on the right are experimentally shown to be the most relevant for prediction

Explanations



“AI predicted that patient X can safely stop treatment with confidence score of 0.8. Why? Can I trust it?”



“AI did a code review on my pull request and rejected it. Why? What should I change to get it merged?”



“AI denied loan to applicant A while approved it for applicant B although their profiles look similar. Why? Can I trust it?”

Regulations

- B. Goodman and S. Flaxman. EU regulations on algorithmic decision-making and a ‘right to explanation’. In Proc. ICML Workshop Human Interp. Mach. Learn., pages 26–30, New York, NY, June 2016
- D. B. Pasternak. Illinois and City of Chicago poised to implement new laws addressing changes in the workplace — signs of things to come? The National Law Review, June 2019
- A. D. Selbst and J. Powles. Meaningful information and the right to explanation. Int. Data Privacy Law, 7(4):233–242, Nov. 2017
- K. R. Varshney. Trustworthy machine learning and artificial intelligence. ACM XRDS Mag., 25(3):26–29, Spring 2019

General Data Protection Regulation (GDPR)

- Articles 13 and 14 state that a data subject has the right to “meaningful information about the logic involved”
- Recital 71 states more clearly that a person who has been subject to automated decision-making “should be subject to suitable safeguards” which should include
 - specific information to the data subject
 - the right to obtain human intervention to express his or her point of view
 - to obtain an explanation of the decision reached after such assessment
 - and to challenge the decision

Explainable AI

- Tools and techniques for humans to
 - Understand rationale behind AI systems' predictions
 - Establish trust in AI systems involved in making decisions
- In a nutshell, we want to open the ML black box and make it *interpretable* other than accurate

Taxonomy

Features VS Samples

Local VS Global

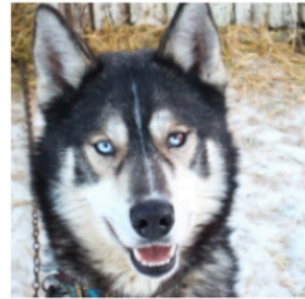
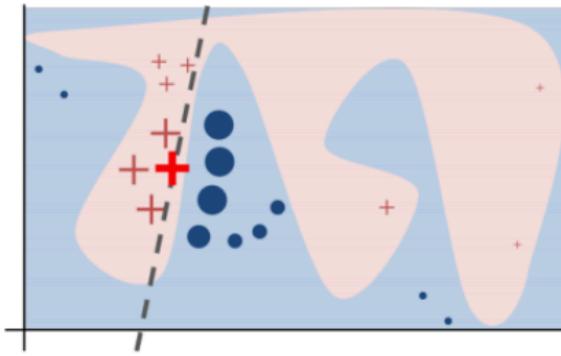
Static VS Interactive

Directly interpretable VS Post-hoc

Surrogate VS Visualisation

Black Box VS White Box

Explanations via features



(a) Husky classified as wolf



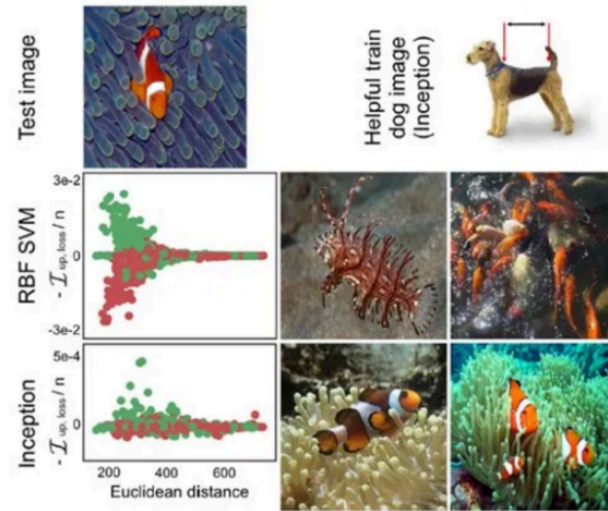
(b) Explanation

- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Why should i trust you?: Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, 2016.

Features VS
Samples

Explanations via samples

- "What would happen if a given training point was not available?"
- "What would happen if we would change a training point values of a small amount?"
- The **influence function** is a measure of how strongly the model parameters or predictions depend on a training instance without retraining the whole model
- Koh, Pang Wei, and Percy Liang. "Understanding black-box predictions via influence functions." Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017.



Features VS
Samples

Local VS Global



Local: For describing the behaviour of a single prediction



Global: For describing the behaviour of the entire model

Binning | X-Axis

Inference sco ▼

Count

10

Binning | Y-Axis

(none) ▼

Color By

COMPASS_d ▼

Label

(d)

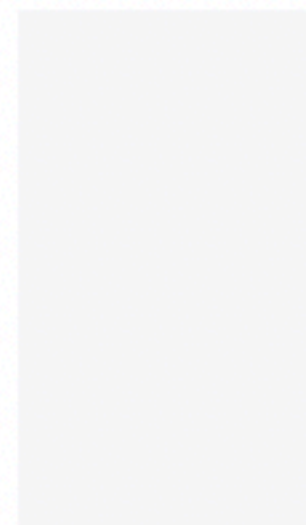
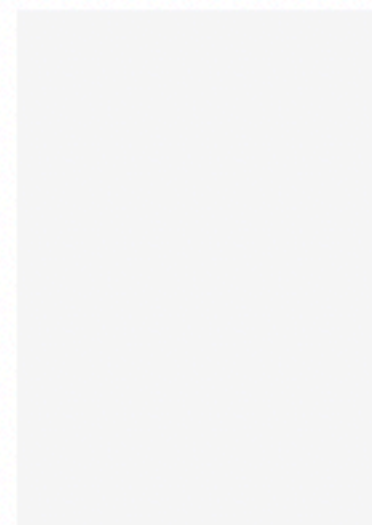
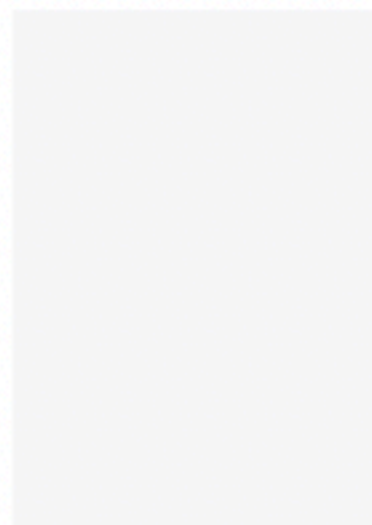
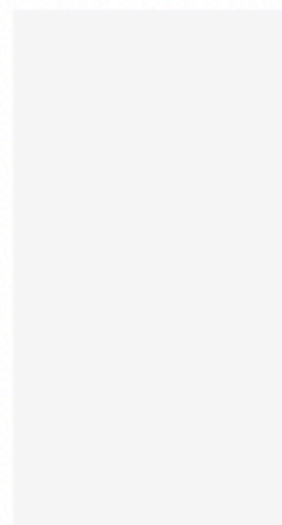
0.0225 — 0.12

0.12 — 0.217

0.217 — 0.315

0.315 — 0.412

0.412 — 0.51



Directly interpretable VS post-hoc



Directly interpretable: By its intrinsic transparent nature the explanation is understandable by most consumers (e.g. a small decision tree)



Post-hoc: The explanation involves an auxiliary method to explain a model after it has been trained



Surrogate. A second, usually directly interpretable, model that approximates a more complex (and less interpretable) one , e.g., a regression model



Visualisation. A focus on parts of a model that are more easily understandable , e.g., deep dream

Surrogate VS Visualisation

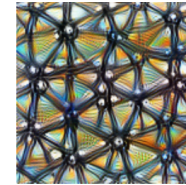
Feature Visualization by Optimization

Different **optimization objectives** show what different parts of a network are looking for.

n layer index
x, y spatial position
z channel index
k class index



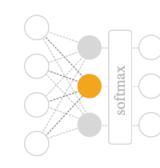
Neuron
 $\text{layer}_n[x, y, z]$



Channel
 $\text{layer}_n[:, :, z]$



Layer/DeepDream
 $\text{layer}_n[:, :, :]^2$

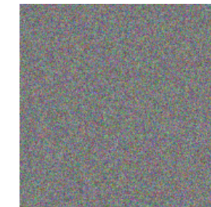


Class Logits
 $\text{pre_softmax}[k]$

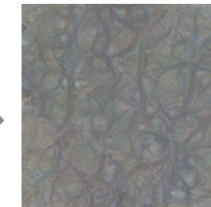


Class Probability
 $\text{softmax}[k]$

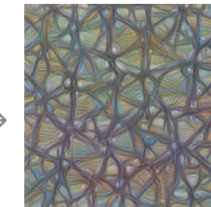
Starting from random noise, we optimize an image to activate a particular neuron (layer mixed4a, unit 11).



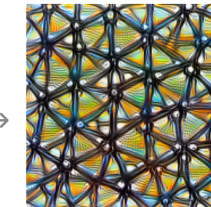
Step 0



Step 4



Step 48



Step 2048

Black Box vs White Box

Black-box methods come with a model-agnostic interface, e.g., by perturbing input data

White-box rely on the internal mechanisms of the model, e.g., by backpropagating the contributions of all neurons in the network to every feature of the input.

The background of the slide features several thin, curved lines in shades of gray, some solid and some dashed, creating a sense of motion or a stylized globe. A blue rectangular box with a speech bubble tail is positioned on the left side.

Other categories

- **Counterfactual explanations**
- **Causal explanations**
- **Explanations aggregators**
- **...**

A new dimension

We introduce task-specific techniques, as opposite to previous (task-agnostic) techniques

Such category is inspired from specific data integration task, where the nature of the problem makes previous techniques ineffective

The background features a series of concentric circles in a light blue-grey color. A dashed line of the same color starts from the left edge and curves around the text, ending towards the bottom right.

▼ Rapid zoom to our
specific task

Data Integration

We aim at providing a unified view over data



We get data from multiple, autonomous, sources

**e.g., in the domain of e-commerce, we have alibaba,
amazon, etc**



We produce a holistic data structure for supporting advanced tasks

e.g., question answering, search,

Data integration pipeline



ENTITY RESOLUTION (ER)

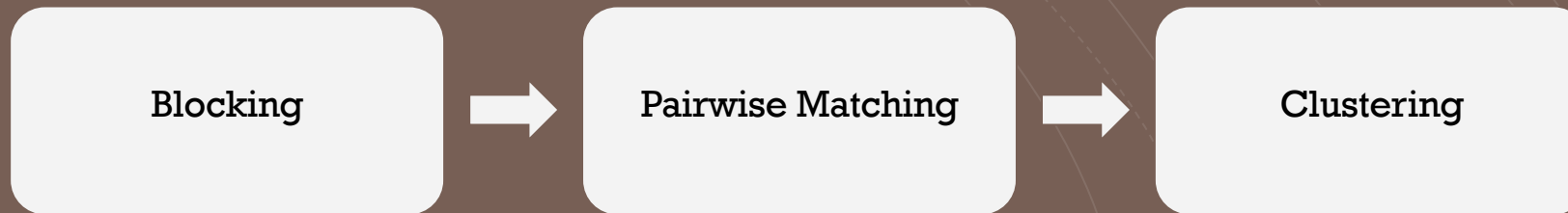
Problem definition:

Consider a set of data sources S , providing a set of records R over a set of attributes A . Entity Resolution computes a partitioning P of R , such that each partition in P identifies the records in R that refer to a distinct entity.

s_1	Products				
	Brand	Model	Resolution	Digital Zoom	Optical Zoom
r_1	Sony	Alpha 7	16mpx	16x	8x
r_2	Sony	ILCE 7	16.0 MP	16x	8x

r_n	Sony	Alpha 5	8.0 MP	8x	4x
	a_1	a_2	a_3	a_4	a_5

Entity resolution pipeline



Short history of ER solutions

growing *F*-Measure and Scale



<i>approach</i>		<i>example techniques</i>		
		Blocking	Pair Matching	Clustering
~1970	Rules & Stats	same name	string similarity	trans. closure
~2000	Supervised / Unsup Learning		decision trees	corr. clust.
~2015	Supervised Learning	active learn.	random forests	
~2018	Deep Learning	embeddings	DNNs	

source: https://thodrek.github.io/di-ml/vldb2018/slides/diml_vldb2018.pdf

Short history of ER solutions

growing **OPACITY**

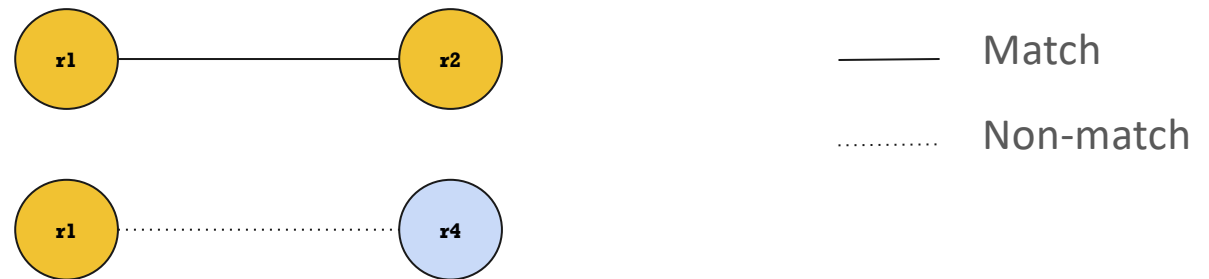


example techniques

		Blocking	Pair Matching	Clustering
~1970	Rules & Stats	same name	string similarity	trans. closure
~2000	Supervised / Unsup Learning		decision trees	corr. clust.
~2015	Supervised Learning	active learn.	random forests	
~2018	Deep Learning	embeddings	DNNs	

ER: PAIRWISE MATCHING

Basic step of ER: compares a *pairs of records* and makes a local decision of whether or not they refer to the same entity.



DEEP LEARNING FOR ENTITY RESOLUTION

- Two main systems:
 - **DeepMatcher**, a modular architecture for record linkage
 - **DeepER** , a specific architecture for record linkage and a blocking system

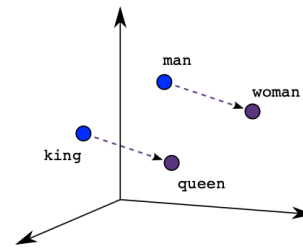
KEY CONCEPTS

Use pre-trained word-embedding models to represent tokens in the dataset, such as Glove, FastText or Word2vec

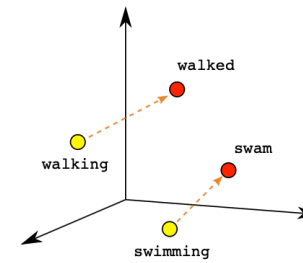
Reuse well known techniques for NLP processing, such as **RNN** or **LSTM**, to summarize attribute tokens

Exploit the ability of deep learning to approximate very complex functions

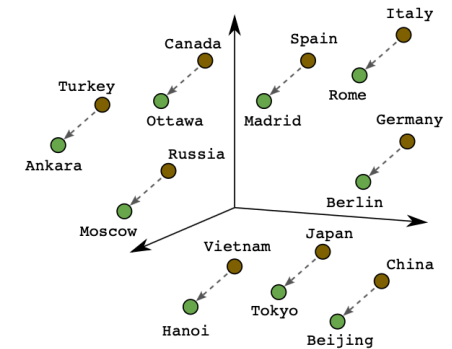
Word Embedding



Male-Female

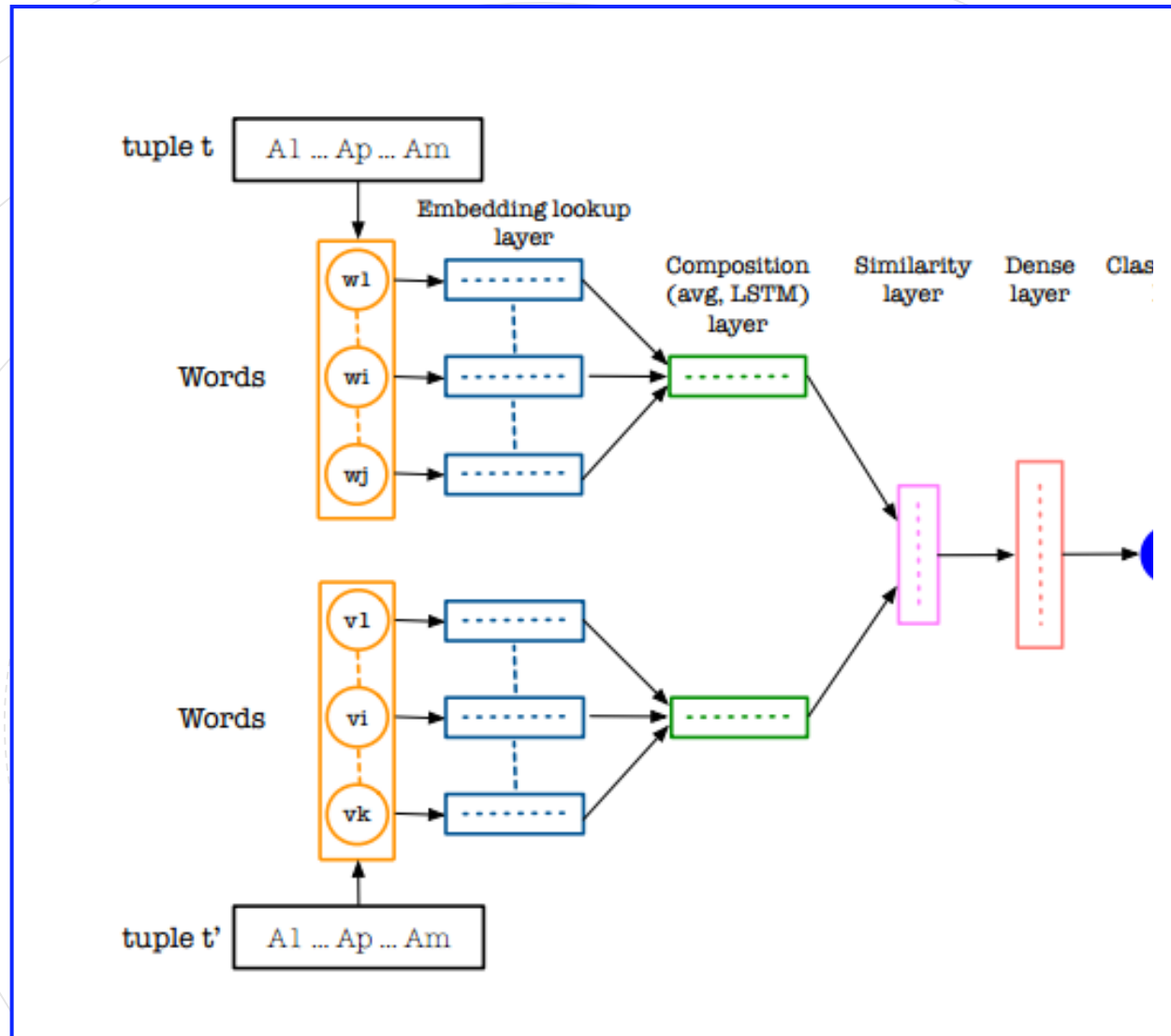


Verb Tense



Country-Capital

- **Collective name for a set of language modeling and feature learning techniques in natural language processing (NLP)**
- **Words or phrases from the vocabulary are mapped to vectors of real numbers, keeping the semantic**



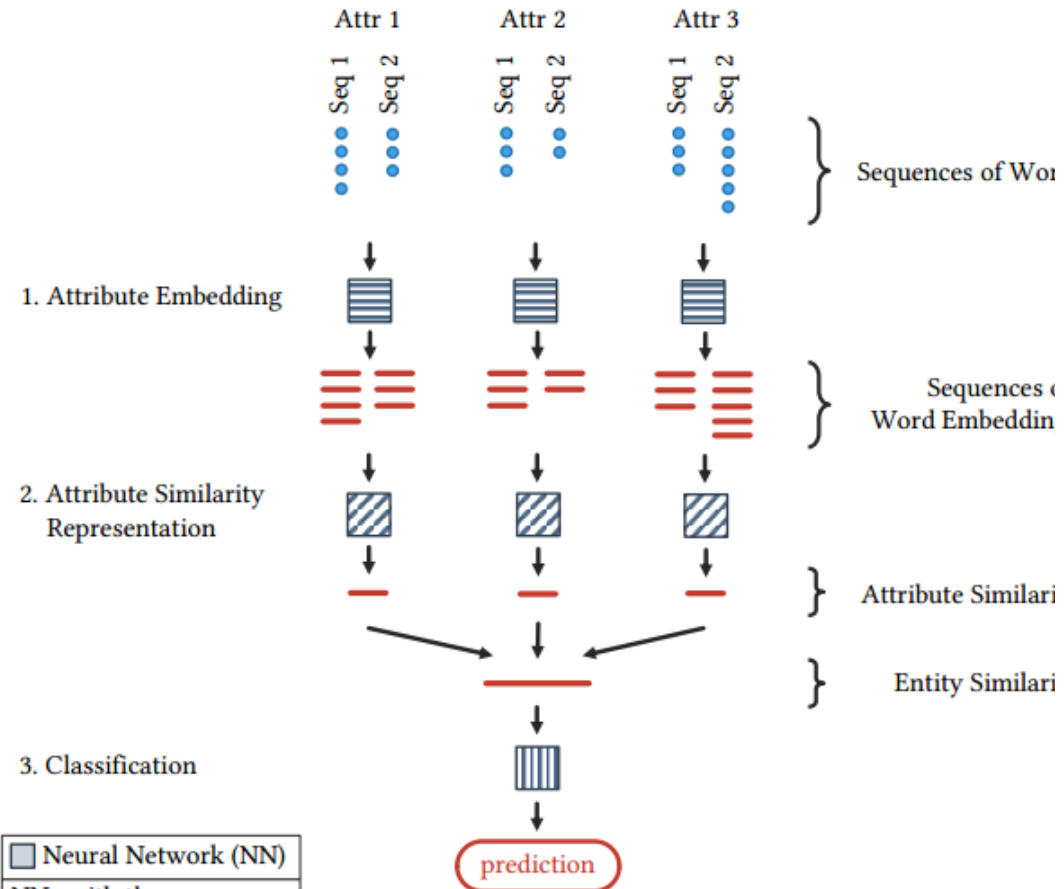
DEEPER
ARCHITECTURE

NEGATIVE SAMPLES BUILDING IN DEEPER

- In DeepER the negative training samples (pair of non-matching records) are built in the following way:
 - Let S_1 and S_2 the two sources of records
 - Take random pairs from S_1, S_2 and evaluate their similarity with common string similarity functions, such as Levenshtein, Jaro-Winkler, Jaccard etc..
 - If the similarity of a pair is under a threshold T the pair is considered a negative sample
- This technique allows to increase dramatically the performances (in terms of F1-score)

BLOCKING FUNCTION

- Once the model is trained, the output of the **LSTM layer** can be used to do blocking
- Let v_1 and v_2 two input records and let v_1' and v_2' the representations of v_1 and v_2 from the LSTM output
- Given a generic hash function h the assumption is that $h(v_1') = h(v_2')$ if v_1 and v_2 refer to the same entity



DEEPMATCHER ARCHITECTURE

(<https://github.com/anhaidgroup/deepmatcher>)

DEEPMATCHER

MAIN FEATURES

- Different component for each attribute of the dataset
- 4 different ways of attribute summarization (SIF, RNN, Attention and Hybrid)
- Custom classification and attribute comparison layers
- Not-trainable embedding layer, with the possibility to choose pre-trained model (FastText or Glove)

Opacity of Pair-Matching Models

94.9% ✓

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
ITunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

model: DeepMatcher

Opacity of Pair-Matching Models



	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

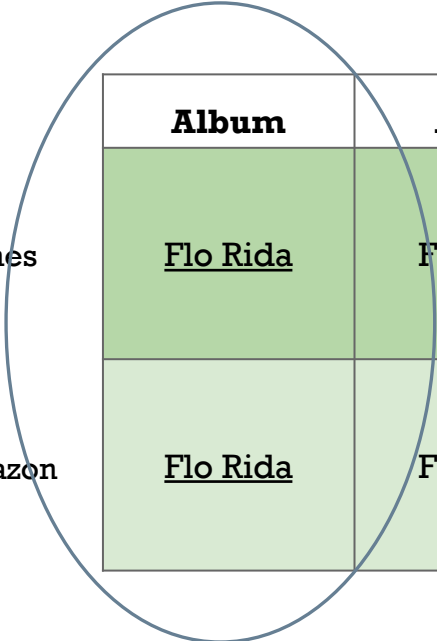
What feature has the model picked?

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	<u>Flo Rida</u>	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	<u>Flo Rida</u>	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

What feature has the model picked?

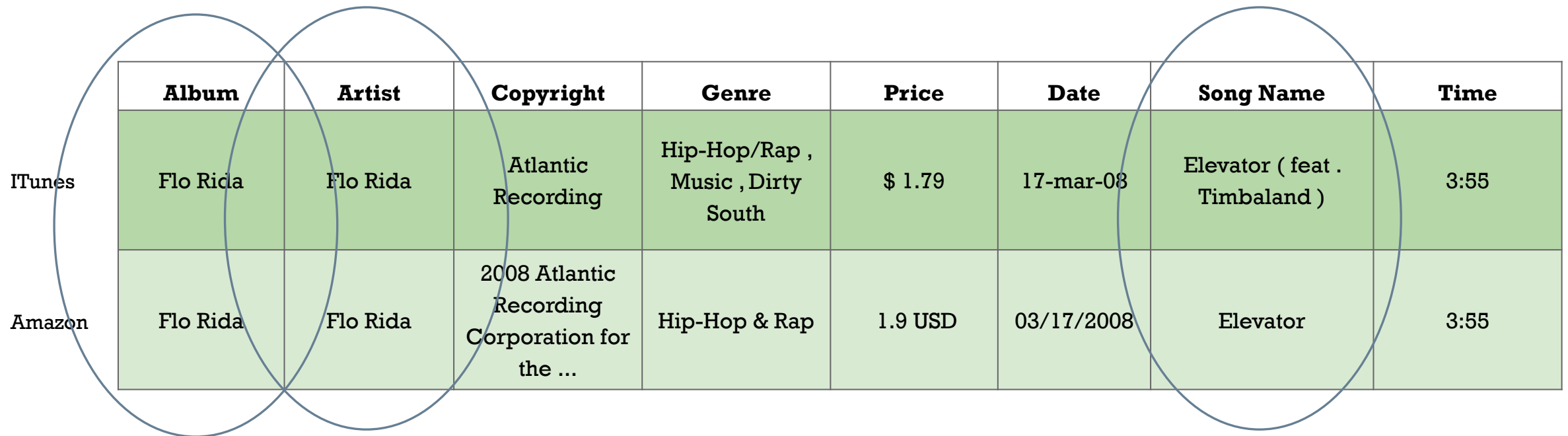
	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	<u>Elevator</u> (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	<u>Elevator</u>	3:55

What feature has the model picked?



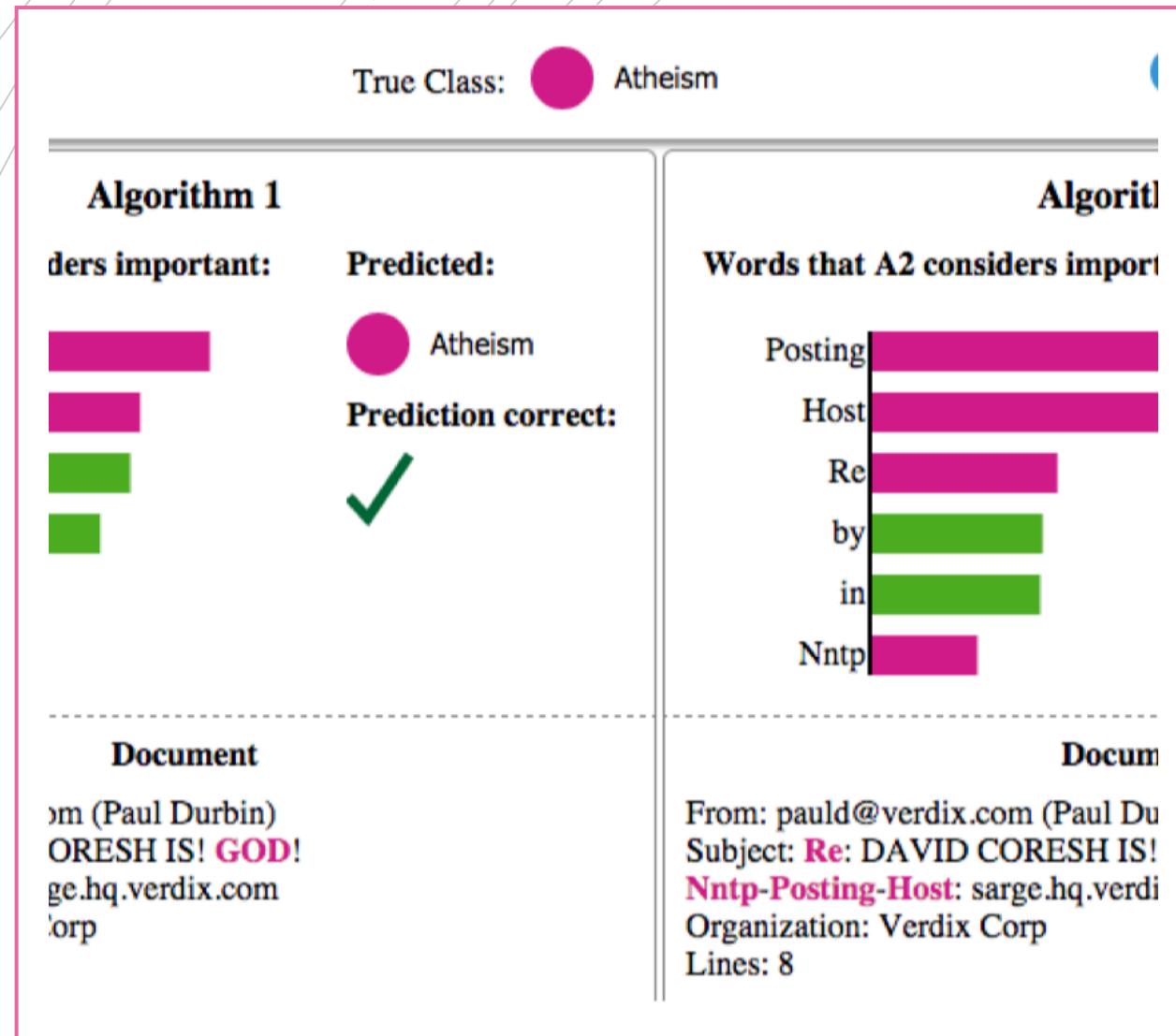
	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	<u>Flo Rida</u>	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	<u>Flo Rida</u>	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

What feature has the model picked?



	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

Popular Model-Agnostic Explanation Tool: LIME



- LIME is a framework for explaining individual predictions of black-box models
- ← “Christianity” or “Atheism”

Given a prediction, it considers an interpretable feature space, e.g. its tokens.

GOD	Mean	Anyone	This	Koresh	through
1	1	1	1	1	1

Class = «Christianity»

It makes random perturbations, for instance by dropping tokens.

GOD	Mean	Anyone	This	Koresh	through
1	0	1	0	1	1

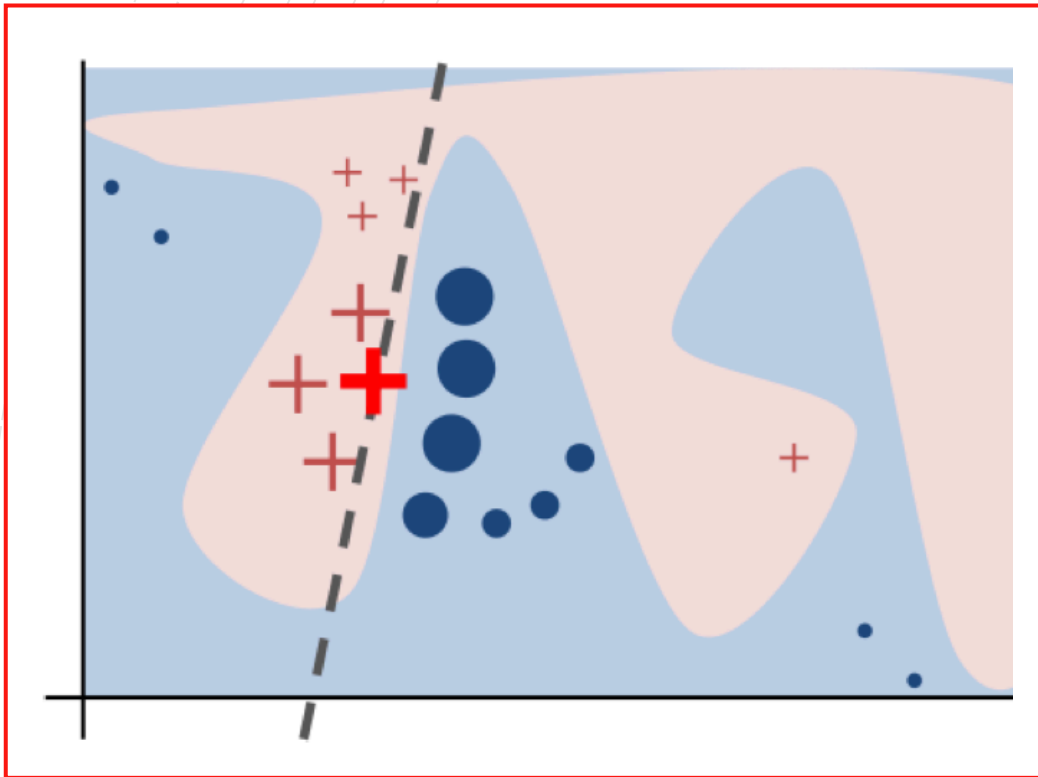
Class = «Christianity»

And tracks how the prediction changes from perturbations, for instance by dropping tokens

GOD	Mean	Anyone	This	Koresh	through
0	1	1	1	1	1

Class = «**Atheist**»

Surrogate model



- Finally, it learns an surrogate interpretable model of the perturbed instances predictions, e.g. a linear regression model
- The interpretable model is not faithful globally, but locally can give accurate influence scores of each feature, e.g., tokens

Relevance scores

GOD	Mean	Anyone	This	Koresh	through
0.3	0.02	0.1	0.14	0.16	0.07

Can we use LIME to explain ER predictions?



	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

Mojito = LIME for ER

- Classification Instance = Pair of Records

ITunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

Mojito = LIME for ER

- Specifically, a bag of the original left and right tokens, with a prefix

iTunes=L	Lalbum_Flo	Lartist_Flo	Lcopyright_Atlantic	Lprice_\$	Ldate_17-	Ltitle_Elevator	Ltime_3:55
	Lalbum_Rida	Lartist_Rida	Lcopyright_Recording	Lprice_1.79	mar-08	(Ltitle_feat . Ltitle_Timbaland)	
Amazon=R	Ralbum_Flo	Rartist_Flo	Rcopyright_2008	Rprice_1.9	Rdate_03/17	Rtitle_Elevator	Rtime_3:55
	Ralbum_Rida	Rartist_Rida	Rcopyright_Atlantic Rcopyright_Recording ..	Rprice_USD	/2008		

Mojito = LIME for ER

- Behind the scenes, the prefix gives us the ability to perform document perturbations that make more sense for the ER task

Lalbum_Flo Lalbum_Rida Lartist_Flo Lartist_Rida Lcopyright_Atlantic Lcopyright_Recording Lprice_\$ Lprice_1.79 Ldate_17-mar-08 Ltitle_Elevator (Ltitle_feat . Ltitle_Timbaland) Ltime_3:55 Ralbum_Flo Ralbum_Rida Rartist_Flo Rartist_Rida Rcopyright_2008 Rcopyright_Atlantic Rcopyright_Recording ... Rprice_1.9 Rprice_USD Rdate_03/17/2008 Rtitle_Elevator Rtime_3:55


Mojito's perturbations primitives

- In addition, Mojito extends LIME with a new set of perturbation primitives
 - A variant of the original DROP primitive
 - A new COPY primitive

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

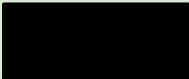
DROP

- The DROP primitive typically DECREASES similarity
- E.g., Remove one token from a matching attribute

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo 	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap	1.9 USD	03/17/2008	Elevator	3:55

DROP

- The DROP primitive typically DECREASES similarity
- E.g., Remove an attribute entirely from one of the records

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap		03/17/2008	Elevator	3:55

Discussion

- The DROP primitive can also INCREASE similarity
- E.g., Remove a non-matching attribute from both

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South		17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap		03/17/2008	Elevator	3:55

COPY

- The COPY primitive always INCREASES similarity
- That is, making two attributes matching or more similar
- Specific for the ER task

	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	Atlantic Recording	Hip-Hop/Rap , Music , Dirty South	\$ 1.79	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop & Rap , Music , Dirty South	1.9 USD	03/17/2008	Elevator	3:55

Putting all together

Mojito considers all the pairs of records in the test set

Applies random DROP/COPY perturbations using the LIME engine

Collects all the influence scores returned by LIME

Returns both

We demonstrate Mojito on two datasets: (1) SONGS and (2) BEERS

aggregate scores of ATTRIBUTE

aggregate scores of TOKEN for each attribute

Time

Song Name

Date

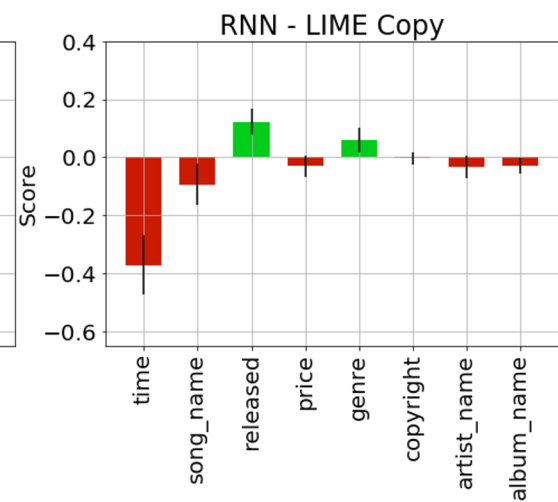
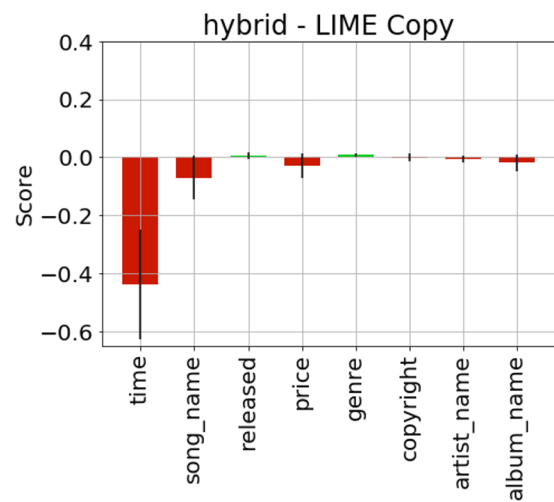
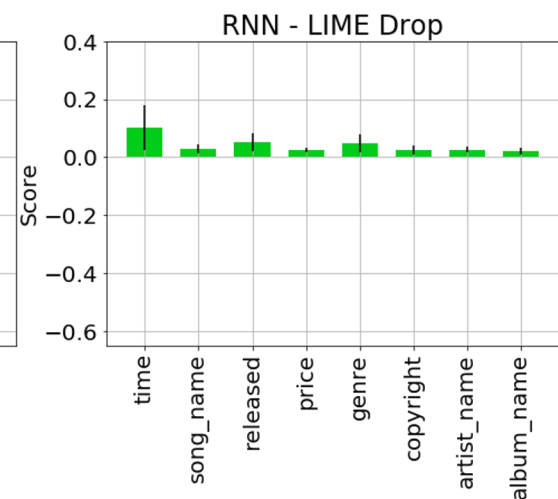
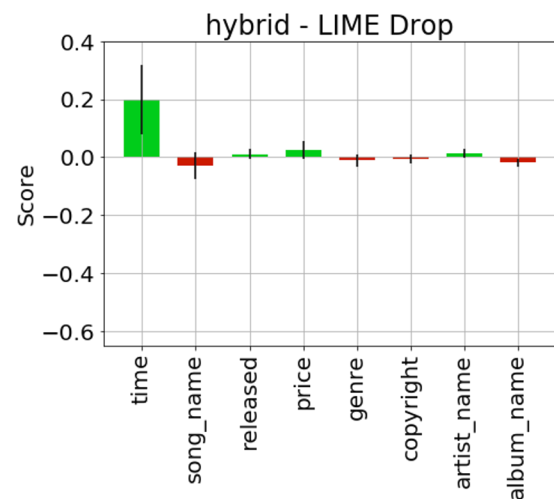
Price

Genre

Copyright

Artist

Album



Manual Check: Non-Match to Match

- Take a non-matching pair



	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
iTunes	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop/Rap , Music , Dirty South	\$ 1.99	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	*	*	*	*	*	*	*	*

Manual Check: Non-Match to Match

- Set TIME to the same (or close) value and it becomes a match



	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
ITunes	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop/Rap , Music , Dirty South	\$ 1.99	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	*	*	*	*	*	*	*	3:55

Observations

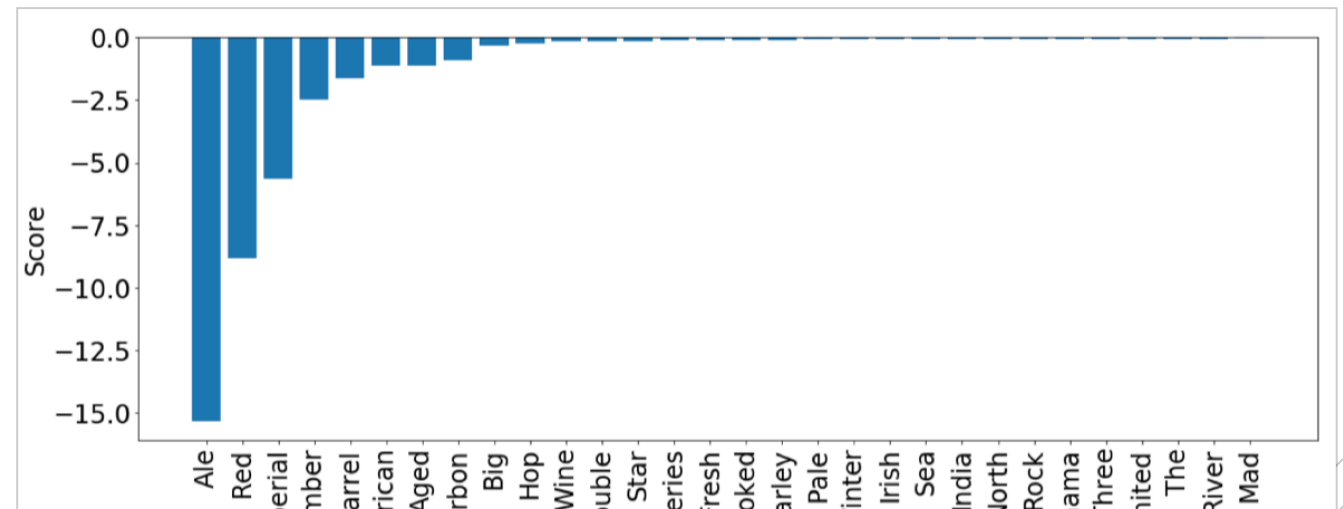
- Most **matching pairs** in the training set have same time



	Album	Artist	Copyright	Genre	Price	Date	Song Name	Time
ITunes	Flo Rida	Flo Rida	2008 Atlantic Recording Corporation for the ...	Hip-Hop/Rap , Music , Dirty South	\$ 1.99	17-mar-08	Elevator (feat . Timbaland)	3:55
Amazon	*	*	*	*	*	*	*	3:55

Mojito's token-level scores for Beer Name

- Imperial red ale on the top



Manual Check: Match to Non-Match

- Take a matching pair



ABV	Beer Name	Brewery	Style
5.60%	Sanibel Red Island Ale	Point Ybel Brewing Company	American Amber / Red Ale
5.60%	Point Ybel Sanibel Red Island Ale	Point Ybel Brewing Company	Irish Ale

Manual Check: Match to Non-Match

- Make “Imperial Red Ale” appear in the Beer Name and it becomes a non-match



ABV	Beer Name	Brewery	Style
5.60%	Sanibel Red Island Imperial Red Ale	Point Ybel Brewing Company	American Amber / Red Ale
5.60%	Point Ybel Sanibel Red Island Imperial Red Ale	Point Ybel Brewing Company	Irish Ale

Manual Check: Non-Match to Match

- Take a non-matching pair involving two Imperial Red Ales



ABV	Beer Name	Brewery	Style
9.00 %	Hop Around Imperial Red Ale	Big Bay Brewing Co.	American Amber / Red Ale
9.00 %	Marble Imperial Red Ale	Marble Brewery	American Strong Ale

Manual Check: Non-Match to Match

- Remove “Imperial Red Ale” from the Beer Name and it becomes a match
- Even though they still look very different



ABV	Beer Name	Brewery	Style
9.00 %	Hop Around Imperial Red Ale	Big Bay Brewing Co.	American Amber / Red Ale
9.00 %	Marble Imperial Red Ale	Marble Brewery	American Strong Ale

Observations

- Most **non-matching pairs** in the training set involve Imperial Red Ales



ABV	Beer Name	Brewery	Style
9.00 %	Hop Around Imperial Red Ale	Big Bay Brewing Co.	American Amber / Red Ale
9.00 %	Marble Imperial Red Ale	Marble Brewery	American Strong Ale

Conclusions

Explainable AI is an exciting field

Many opaque data integration models that need to be equipped with explainable tools

Mojito is an extension of LIME for the specific ER tas

It builds on two main intuitions

- represents pairs of records as a single document, in order to leverage the LIME framework
- plugs in ER specific perturbations

Explanations can be used to «debug» the model

Project

Run

Run NLP based models for ER (e.g., based on BERT) over our benchmark (more on next slide)

Run

Run traditional Machine Learning models for ER (e.g., Magellan)

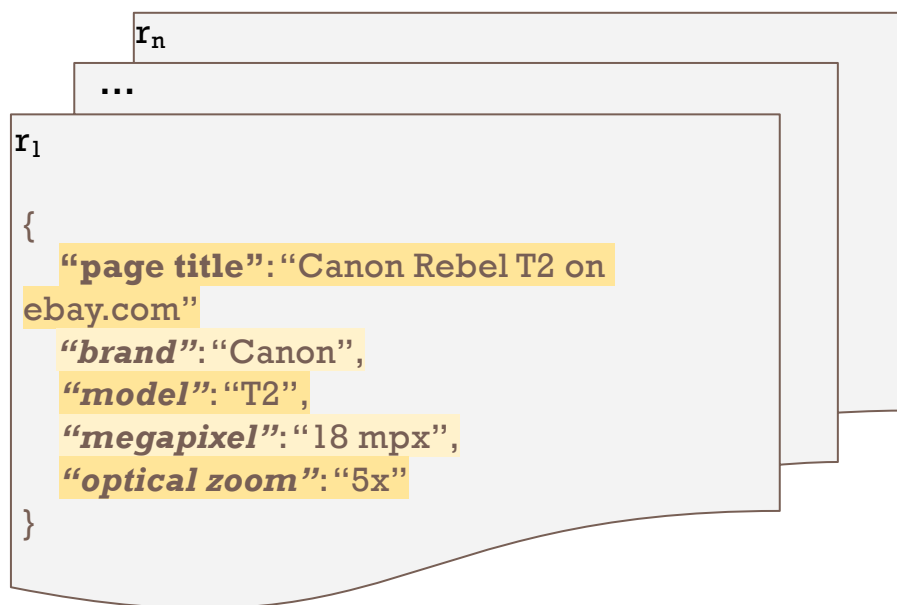
Find

Find unexpected behaviors and try to explain them

ALASKA BENCHMARK



End-to-end benchmark for Big Data Integration tasks based on *real-world product specification*



Available datasets:

Dataset	# data sources	# records	# distinct attributes
<i>CAMERA</i>	24	~30k	~4k
<i>MONITOR</i>	26	~16k	~2k

<http://alaska.inf.uniroma3.it/>

Credits



Tommaso Teofili
PhD Student
(Explanation taxonomy)



Vincenzo Martello
(ER models)



Andrea De Angelis,
Research
(DI pipeline & Alaska)



Vincenzo di Cicco
(Mojito)

References

1. Doshi-Velez, Finale, and Been Kim. "Towards a rigorous science of interpretable machine learning." arXiv preprint arXiv:1702.08608 (2017)
2. Arya, Vijay, et al. "One Explanation Does Not Fit All: A Toolkit and Taxonomy of AI Explainability Techniques." arXiv preprint arXiv:1909.03012 (2019).
3. Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "" Why should i trust you?" Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016.
4. Koh, Pang Wei, and Percy Liang. "Understanding black-box predictions via influence functions." *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017.
5. Wexler, James, et al. "The What-If Tool: Interactive Probing of Machine Learning Models." *IEEE transactions on visualization and computer graphics* (2019).
6. Dong, Xin Luna, and Theodoros Rekatsinas. "Data integration and machine learning: A natural synergy." *Proceedings of the 2018 International Conference on Management of Data*. 2018.
7. Ebraheem, Muhammad, et al. "DeepER--Deep Entity Resolution." *arXiv preprint arXiv:1710.00597* (2017).
8. Mudgal, S., Li, H., Rekatsinas, T., Doan, A., Park, Y., Krishnan, G., ... & Raghavendra, V. (2018, May). Deep learning for entity matching: A design space exploration. In *Proceedings of the 2018 International Conference on Management of Data* (pp. 19-34).
9. Di Cicco, Vincenzo, et al. "Interpreting deep learning models for entity resolution: an experience report using LIME." *Proceedings of the Second International Workshop on Exploiting Artificial Intelligence Techniques for Data Management*. 2019.
10. Gryz, Jarek, and Nima Shahbazi. "Futility of a Right to Explanation." *PIE Workshop @ EDBT/ICDT 2020*



Thanks for your attention