## XAI

### Visual Analytics for explainable Deep Learning

• Learn from data without explicit programming



- Learn from data without explicit programming
- ML algorithms improve performance through experience



- Learn from data without explicit programming
- ML algorithms improve performance through experience
  - Learn from data and generalize to unseen data



- Learn from data without explicit programming
- ML algorithms improve performance through experience
  - Learn from data and generalize to unseen data
- Key goal: Create systems that can automatically learn patterns and make intelligent decisions



## **Types of Machine Learning**

- Supervised Learning: Learns from labelled training data
  - Examples: Classification, Regression
  - Used in spam detection, price prediction
- Unsupervised Learning: Finds patterns in unlabelled data
  - Examples: Clustering, Dimensionality Reduction
  - Used in customer segmentation, anomaly detection
- Reinforcement Learning: Learns through trial and error
  - Used in game AI, robotics, autonomous systems

## **Deep Learning**

- Advanced machine learning technique using multi-layered neural networks
- Mimics human brain's neural structure to process complex data
- Enables automatic feature extraction and learning from raw data

• Layers: Input, Hidden, Output layers



- Layers: Input, Hidden, Output layers
- Neurons: Process and transmit information



- Layers: Input, Hidden, Output layers
- Neurons: Process and transmit information
- Activation Functions: Introduce non-linearity

	Input Layer
Reurons	
Activation	Be Hidden Layer
Functions	은 Output Layer

- Layers: Input, Hidden, Output layers
- Neurons: Process and transmit information
- Activation Functions: Introduce non-linearity
- Backpropagation: Learns by adjusting network weights



- Convolutional Neural Networks (CNNs)
  - Image and video processing
  - Pattern recognition



U-Net: Convolutional Networks for Biomedical Image Segmentation, https://doi.org/10.48550/arXiv.1505.04597

- Recurrent Neural Networks (RNNs)
  - Sequential data analysis
  - Natural language processing



Feedforward Neural Network

Recurrent Neural Network

Recurrent Neural Networks (RNNs): A gentle Introduction and Overview, https://doi.org/10.48550/arXiv.1912.05911

- Transformers
  - Advanced language models
  - Context understanding



Attention Is All You Need, https://doi.org/10.48550/arXiv.1706.03762

## **Training Deep Learning Models**

- Large datasets critical for performance
- Requires significant computational power
- Techniques: Transfer Learning, Data Augmentation
- Challenges: Overfitting, Computational Complexity

## Why We Need Explainable AI

- It is difficult to trust a system whose decisions we don't understand.
- Lack of transparency can lead to bias and unfair outcomes.
- Explainability is essential for accountability and regulation of AI.

### Learning Performance vs. Explainability



### Overview



Figure 1. Overview of explainable deep learning

[Choo]

.....

## Definitions

- XAI aims to develop machine learning techniques that provide understandable, trustworthy and explainable rationales for decisions made by black-box models. [Adadi, Dragoni, Gunning]
- XAI is the class of systems that provide visibility into how an AI system makes decisions and predictions and executes its actions. [Rai]
- XAI refers to systems that try to explain how a black-box AI model produces its outcomes. [Moradi]

## Definitions

- XAI aims to develop machine learning techniques that provide understandable, trustworthy and **explainable** rationales for decisions made by black-box models. [Adadi, Dragoni, Gunning]
- XAI is the class of systems that provide visibility into how an AI system makes decisions and predictions and executes its actions. [Rai]
- XAI refers to systems that try to **explain** how a black-box AI model produces its outcomes. [Moradi]

-> **Explain** how Black-Box-Models get to their decisions

## What is XAI

- XAI refers to techniques that make AI systems more understandable to humans.
- It aims to provide insights into how AI systems make decisions and why they reach conclusions.

### Stakeholders / Who needs XAI?

- Data scientists use XAI to debug and improve models.
- **Domain experts** need to understand the reasoning behind AI's recommendations.
- **Decision-makers** rely on XAI to ensure responsible use of AI.
- End users deserve to know how AI affects their lives.

### **Important Questions**

#### §4 WHY

Why would one want to use visualization in deep learning?

#### **§6 WHAT**

What data, features, and relationships in deep learning can be visualized?

#### §8 WHEN

When in the deep learning process is visualization used?



#### [Homann]

### **Important Questions**

#### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

#### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

#### §8 WHEN

When in the deep learning process is visualization used?



§5 WHO

Who would use and benefit from visualizing deep learning?

#### 7 HOW

How can we visualize deep learning data, features, and relationships?



Where has deep learning visualization been used?

#### [Homann]

### **Important Questions**

#### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

#### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

#### §8 WHEN

When in the deep learning process is visualization used?





[Homann]

Who would use and benefit from visualizing deep learning?

Model Developers & Builders Model Users Non-experts



How can we visualize deep learning data, features, and relationships?



Where has deep learning visualization been used?

Fig. 1 An overview of visual analytics research for machine learning.



Fig. 1 An overview of visual analytics research for machine learning.

Data Preparation Feature Extraction

Fig. 1 An overview of visual analytics research for machine learning.



Fig. 1 An overview of visual analytics research for machine learning.



Fig. 1 An overview of visual analytics research for machine learning.



Fig. 1 An overview of visual analytics research for machine learning.



Fig. 1 An overview of visual analytics research for machine learning.








# Overview – Visual Analytics for Machine Learning Machine Learning Pipeline





# Overview – Visual Analytics for Machine Learning



Fig. 1 An overview of visual analytics research for machine learning.

# Overview – Visual Analytics for Machine Learning



Fig. 1 An overview of visual analytics research for machine learning.

## Overview – Visual Analytics for Machine Learning



Fig. 1 An overview of visual analytics research for machine learning.

## Overview



Fig. 1 An overview of visual analytics research for machine learning.

## Techniques Before Model Building

Ensure high-quality **data** and **features** to improve model performance and reliability.

- Improving Data Quality
- Improving Feature Quality

## Improving Data Quality - Instance-Level

- Anomaly Detection and Correction:
  - Visualising and interacting with data to identify missing values,
  - outliers,
  - duplicates,
  - and out-of-distribution samples.
  - i.e. Profiler and OoDAnalyzer
- Provenance Tracking:
  - Illustrating the impact of data cleaning and preprocessing steps on data quality.
  - i.e. DQProv Explorer
- Privacy Preservation:
  - Balancing data utility with privacy concerns during the cleaning process
  - i.e. Privacy Exposure Risk Tree and GraphProtector.

## Improving Data Quality - Instance-Level

• Anomaly Detection and Correction Example: Profiler



Profiler: integrated statistical analysis and visualization for data quality assessment, https://doi.org/10.1145/2254556.2254659



Source (Levenshtein)

## Improving Data Quality - Instance-Level

• Provenance Tracking Example: DQProvExplorer



Capturing and Visualizing Provenance From DataWrangling, doi: 10.1109/MCG.2019.2941856.







In the Provenance Graph View, the operation icons show that in the orange wrangling branch text transformations were used, opposed to the blue branch, where rows were removed.

Inspection of the heights of both branches' end nodes (see highlighted areas) also shows that the orange branch contains more entries/rows, which means more information has been retained.

## Improving Data Quality - Instance-Level

• Privacy Preservation Example: GraphProtector



GraphProtector: A Visual Interface for Employing and AssessingMultiple Privacy Preserving Graph Algorithms, doi: 10.1109/TVCG.2018.2865021.



## Improving Data Quality - Label-Level

- Noisy Label Handling:
  - Visualising and refining crowdsourced annotations,
  - identifying unreliable workers,
  - correcting mislabelled instances.
  - i.e. LabelInspect and C2A
- Interactive Labelling:
  - Efficiently labelling unlabelled data by leveraging techniques like clustering similar instances and filtering to find items of interest.
  - i.e. MediaTable and the SOM-based visualisation by Moehrmann et al..

## Improving Data Quality - Label-Level

• Noisy Label Handling Example: LabelInspect



An Interactive Method to Improve Crowdsourced Annotations, doi: 10.1109/TVCG.2018.2864843.



## Improving Feature Quality

- Feature Selection:
  - Select useful features that contribute most to the prediction.
  - i.e. DimStiller and SmartStripes
- Feature Construction:
  - Guide the creation of new, more discriminative features.
  - i.e. FeatureInsight

## Improving Feature Quality

### • Feature Selection Example: DimStiller



DimStiller: Workflows for Dimensional Analysis and Reduction, doi: 10.1109/VAST.2010.5652392.









## Improving Feature Quality

• Feature Construction Example: FeatureInsight

Cycling – 🗆 🔂				
Create Name	a new feature	Teach the classifier to recognize	Cycling Pages 60% Correct (+0.0%) Non-Cycling Pages 75% Correct (+1.7%)	69% Accuracy +1.0%
Words	shopping, account, cart, products, checkout, shop	Terms from pages mistaken for Cycling		Terms from confusingly similar Cycling pages
motore	cycles Eat ×	accesso	brands cyclescheme allez checkout	ride
bike ra race, race	ce East × s, bike race, bike tour, tours	= brands Example X	specialized bookmark featured items products repairs basket shop workshop	cycling
cycle, bicy	rde, bisyster, trisyster, bile	service	terms offers clothing privacy triathlon cart account shopping advanced	<ul> <li>bicycling</li> <li>funded</li> </ul>
	A	i workshop i stock	sitemap offering policy bikes categories	<ul> <li>ridden</li> <li>links</li> </ul>
		i dealer	builders colours service bikeshop siwis	membe
		i price i gift	conditions flow vision stock bicycles custom serving womens lazer cycles	<ul> <li>members</li> <li>awesome</li> </ul>
		■ specials	bike	reports
		= bikes		* contacts
		account		† announ

FeatureInsight: Visual Support forError-Driven Feature Ideation in Text Classification, doi: 10.1109/VAST.2015.7347637.



## Overview – Visual Analytics for Machine Learning



Fig. 1 An overview of visual analytics research for machine learning.

## Techniques **During** Model Building

Gain deeper understanding of model workings, diagnose training issues, and steer model behaviour towards desired outcomes.

- Model Understanding
- Model Diagnosis
- Model Steering

## Model Understanding - Parameter Effects

- Understanding Parameter Effects:
  - Visualising how model outputs change with variations in parameter settings.
  - i.e. BirdVis

## Model Understanding - Parameter Effects

### • Understanding Parameter Effects Example: BirdVis



BirdVis: Visualizing and Understanding Bird Populations, doi: 10.1109/TVCG.2011.176.





- Network-Centric Methods:
  - Exploring model structure,
  - visualising neuron activations,
  - and interpreting how different parts of the model contribute to the final output.
  - i.e. CNNVis for convolutional neural networks and LSTMVis for recurrent neural networks.
- Instance-Centric Methods:
  - Analysing individual instances and their relationships,
  - visualising the representation space learned by the model.
  - i.e. Rauber et al.
- Hybrid Methods:
  - Combining network-centric and instance-centric approaches
  - i.e. Summit and ActiVis
- Surrogate Model Explanations:
  - Using simpler, interpretable models to explain the behaviour of more complex models.
  - i.e. RuleMatrix and DeepVID

• Network-Centric Methods Example: CNNVis



Towards Better Analysis of Deep Convolutional Neural Networks, https://doi.org/10.48550/arXiv.1604.07043



• Instance-Centric Methods Example: Rauber et al.



Visualizing the Hidden Activity of Artificial Neural Networks, doi: 10.1109/TVCG.2016.2598838.



### • Hybrid Methods Example: Summit



SUMMIT: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations, https://doi.org/10.48550/arXiv.1904.02323.



• Surrogate Model Explanations Example: DeepVID



DeepVID: Deep Visual Interpretation and Diagnosis for Image Classifiers via Knowledge Distillation, doi: 10.1109/TVCG.2019.2903943.

Teacher View	index Add Selected: 351 Model Accurac	cy: 98.57%
a)	4 $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$ $4$	7 6 6 4 1 5 6 6
	$\begin{array}{c} \mathbf{d} \ $	4 3 1
VAE View D Reset P	lyline Reset Band band width: 15 samples: 512 Generate & Train Student epoch: 10 batch: 128 temp	p: 10
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	d7       d8       d9       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4       4	.118e-7 .691e-6 .616e-6 .576e-4 .137e-1 .912e-5 .570e-7 .231e-3 .417e-4 .845e-1

## Model Diagnosis

- Analysing Training Results:
  - Diagnose issues by visualising classifier performance,
  - identifying fairness issues,
  - exploring potential model vulnerabilities.
  - i.e. Squares and FairSight
- Analysing Training Dynamics:
  - Monitor the training process,
  - detecting anomalies,
  - understand the evolution of model behaviour over time.
  - i.e. DGMTracker and DQNViz
#### Model Diagnosis



Squares: Supporting Interactive Performance Analysis for Multiclass Classifiers, doi: 10.1109/TVCG.2016.2598828





#### Model Diagnosis

• Analysing Training Dynamics Example: DGMTracker



Analyzing the Training Processes of Deep Generative Models, doi: 10.1109/TVCG.2017.2744938.



## **Model Steering**

- Model Refinement with Human Knowledge:
  - Allowing users to interactively refine models by editing prototypes,
  - adding constraints,
  - correcting outputs.
  - i.e. ProtoSteer and ReVision
- Model Selection from Ensembles:
  - Comparing and selecting the best model from a set of candidate models.
  - i.e. BEAMES and RegressionExplorer

#### **Model Steering**

• Model Refinement with Human Knowledge Example: ProtoSteer

A Overview B Negative B Positive	Ra 🗃 Control Panel
Prototypes  Weight  Prediction	Andel Info     Dataset: yelp     Partition: train
32       N       hours to       get       a       office       and       range! (did int care about the poor service         37       worst bloody marys in       the airport       pass on       by       and	#Prototypes: 70 (Avg. Length: 18.2) accuracy: 0.983
1 prefer other toxicos than going to this one due to there thereitie customer servee do in the second to	History C
Someer was         dreat         i         protect         product         product <thproduct< th="">         product         <thprodu< td=""><td><ul> <li>✓ 10 prototypes </li> <li>✓ -7 &gt; 63 </li> </ul></td></thprodu<></thproduct<>	<ul> <li>✓ 10 prototypes </li> <li>✓ -7 &gt; 63 </li> </ul>
23 this place was plessed service was fast and vrimter super zeeozy good sluff	O 63 prototypes O D -1 El
3 chteb food the langues is amating service in a beer list much needed in ditteb t	O 62 prototypes
25 best bbq and service i have had in swhile poriors are huge seylars fresh and hot spread 61 P ? 2 0 F feed with \$\$	
19 best fresh great service and ambience fore this place i highly account this service place for an event meal at a very exercise price i freedul staff and a great furning i	- Contraction
	Con- Marian I
V 61 great food amating service and one hell of a beer list much needed in others !	needed
B1 great food and metors staff make sure you bring extra time with you this place is hopping t	Allower of the second s
great servce great food decently priced menu . try the bacon double deadly !	le fantikeenderen
V 43 omelet was great ! ! ! ! factored and bacon horrible . OTHER with cheese on point .	
B2 this place closes early on exciting i very diagonal make sure to call before arrival this place blows i dirty loud facebulk reducing Mutureli even the ice facebulk was nit woking	
Joxe blues service okay but best OTHER ever I I hated frog legs very bland drinks average	Encomposition .

ProtoSteer: Steering Deep Sequence Model with Prototypes, doi: 10.1109/TVCG.2019.2934267.



#### **Model Steering**

• Model Selection from Ensembles Example: BEAMES



BEAMES: Interactive Multimodel Steering, Selection, and Inspection for Regression Tasks, doi: 10.1109/MCG.2019.2922592.



## Overview – Visual Analytics for Machine Learning



Fig. 1 An overview of visual analytics research for machine learning.

#### [Yuan]

## Techniques **After** Model Building

Help users make sense of model outputs, gain insights from data analysis results, and evaluate model performance in real-world contexts.

- Understanding Static Data Analysis Results
- Understanding Dynamic Data Analysis Results

- Textual Data Analysis:
  - Visualizing topics, clusters, and relationships extracted from text data
  - used techniques:
    - topic modelling
    - word embedding
  - i.e. TopicPanorama , DemographicVis , and cite2vec.
- Other Data Analysis:
  - Extending visual analytics to other data types,
  - E.g., flow fields and multi-dimensional data,
  - used techniques:
    - subspace analysis
    - pattern matching
  - i.e. SMARTexplore

• Textual Data Analysis : DemographicVis



DemographicVis: Analyzing demographic information based on user generated content, doi: 10.1109/VAST.2015.7347631.



#### choose primary axis: Sex Graphs normalization: no

C								and the second second			and the second second
	1	-	11 m Transfe	* <sup>2</sup>	1.70	4.7	Lan Party of	and the second last	a series	No. Inc. Contraction	A. Harrison C.
Linguistic A	23.15%	33.55%	3.01%	6.54%	4.50%	0%	0%	5.35%	0%	4.75%	10.63%
	anx	Dash	Comma	tter	Adure .	200	achieve	SemiC	quant	present	relativ
	verb	Excism	leisure	health	swear	filler	WC	adverb	poseno	negemo	Sme
	Dash	family	ingest	Dash	shehe	body	leisure	evol	AIPut	humans	18er
	Exclam	money	ouert.	work	ROW	they	500	Sine	assert	you .	549
	cognech	relig	WC	nont	SemiC	posemo	verb		money	leisure	affect
Topic	8.54%	17.83%	13.40%	9.75%	9.83%	0%	12.60%	0%	0%	0%	14.37%
Subreddt	32.62%	32.49%	29.78%	13.04%	0%	21.63%	13.01%	0%	3.31%	0.25%	25.09%
Otherinto	9.14%	3.15%	0%	1.47%	1.09%	12.28%	7.57%	0%	2.84%	3.53%	13.66%
Linguistic	80.03%	05.22%	49.50%	49,20%	55.52%	48,72%	54,30%	50.78%	29.04%	66.27%	62.86%
Topic	80.58%	85.90%	49.31%	50.61%	\$7.65%	49.02%	52.26%	10.71%	41.04%	65.34%	63.75%
Subreadil	82.94%	90.74%	63.85%	62.09%	57.05%	57.50%	67.37%	50.37%	\$7.42%	64.36%	61.61%
Otherinfo	83.57%	89.81%	65.08%	61.90%	56.12%	56.52%	68.32%	50.00%	53.54%	64,25%	66.70%

User	- 0	Title	ų.	Seitheat	i.	Topic 0	4		downs	4	Time	ġ.	ì
chee		Rancho Relaxo is still broken		Sgt.Vou must have a crappy video card. Stuck teurts a		12		4	0		9/0/2014		1
0.0***		Got a new refurbished lightigs, had to test the video car		Don't get me wrong, HL2 is great, but if you're really		Haff, Re		P	0		2/5/2014		
a		I think I found an Easter egg picture of the dev?		I heard that that was due to a gillch with the video ca		outside		7	0		0/4/2014		
*2***		[Steam] Community Choice Wenner: Bastion @ 53.74 (75% o		My old leptop had trouble running Minecraft, but could		GameDeals		6	0		11/0/2012		
801 <sup>-00</sup>		First clean setup with a large window. What do you thin		It's not as bed as you might their, I had some bad fit		battlestations		5	٥		8/3/2014		
-		Looking to build a Gaming PC for about \$3,000 US		Please read the summary below! Also, note that I have n		buildapcforme		5	0		8/2/2013		
		(Build Help) First time PC builder, looking for feedbac		(PCPartPicker part Isig).htp://pcpartpicker.com/p/1Ggk		buildape		5	0		8/5/2013		
10.000		(Build Ready) My ested gaming computer build, \$900-100		EDIT: Allight, here's my re-edited build. (Also, I'm no		buildape		4	0		4/0/2013		
		5500 PC for the Family		(PCPartPicker part las) http://pcpartpicker.com/p/10/JML_		buildapoforme		3	0		9/2/2013		
ei		(Build Help) Looking to upgrade \/deo Card/CPU		Upgrading the video card will give you much more bang f		buildape		3	0		6/6/2013		
ni		(Build Ready)Gaming PC for ~\$1000, first "graft"		Non-overslocking (PCPer/Picker part leg)/tgu/lpcp		buildape		3	0		8/40013		

Totle: First clean setup with a large window. What do you think reddit? Plasti dipped my h100i to be white.

Soldney: I's not as had as you might think. Thad some bad fittings when I pot my loop together, and I ended up with -15ml of water on the PCB of my while the computer was running. I just scaled up the water, replaced the fittings with good case, and it works just fine I year later. You just have to be use conductive liquids.

- Offline Analysis: Exploring patterns and trends in data over time, with all data available before analysis.
  - **Topic Analysis**: Visualizing topic evolution using techniques, which often employing a river metaphor to convey changes over time.
    - i.e. ThemeRiver and TextFlow
  - **Event Analysis**: Revealing important sequential patterns in event data, which leverage techniques like tensor decomposition and stage analysis.
    - i.e. EventThread
  - **Trajectory Analysis**: Visualizing and understanding movement patterns, often using techniques like clustering, pattern mining, and semantic enrichment.
    - i.e. Kruger et al. and Chen et al. .

• Offline Analysis Example: TextFlow



TextFlow: Towards Better Understanding of Evolving Topics in Text, doi: 10.1109/TVCG.2011.239.



- Online Analysis:
  - Tackling streaming data where new data arrives continuously.
  - Area presents challenges for visualizing evolving patterns and integrating with real-time analysis algorithms.
    - i.e. TopicStream

• Online Analysis Example: TopicStream



Online Visual Analytics of Text Streams, doi: 10.1109/TVCG.2015.2509990.



# Visual-based XAI (vXAI)

- Model usage
  - Feature-Based Methods
  - Rule-Based Methods
  - Propagation-Based Methods
  - Case-Based Methods
- Visual Approaches
  - Data representation
  - Local Explanations
  - Global Explanations

#### **Feature-Based Methods**

#### Identifying the Key Factors Driving AI's Decisions

- Methods identifying which input features were most influential in the model's decision.
- Methods like LIME and SHAP are commonly used for featurebased explanations.

#### **Rule-Based Methods**

#### **Explaining Predictions with IF-THEN Rules**

- Techniques expressing the model's logic as a set of humanreadable rules.
- Bayesian Rule Lists (BRL) are often used to generate rule-based explanations.

#### **Propagation-Based Methods**

#### **Tracing Information Flow in Neural Networks**

- Techniques analysing how information flows through a network's layers to identify important features.
- Methods like Saliency Maps, LRP, and Integrated Gradients are used for propagation-based explanations.

#### **Case-Based Methods**

#### Learning from Similar Past Examples

- Methods finding past cases that are like the current input and show how the model treated those cases.
- Case-Based Reasoning (CBR) is a common approach for this type of explanation.

#### Data representation

- Known visualization techniques for data representation
  - Sankey Diagrams
  - Scatterplots
  - Table based visualizations
  - Etc. -> see Visualization lecture and previous lectures in Visual Analytics

#### Local vs. Global Explanations

#### Explaining Individual Predictions vs. Model Behaviour

- **Local explanations**: focus on understanding a specific prediction.
- **Global explanations**: aim to provide insights into the overall model logic and behaviour.

## XAI for Responsible AI

- XAI is crucial for ensuring that AI is developed and used responsibly.
- It promotes transparency, fairness, accountability, and ethical considerations in AI.

# LIME - Local Interpretable Model-Agnostic **E**xplanation

- Explains predictions of any classifier
- Learn interpretable model locally around the prediction
- Model understanding depends on
  - 1. Trusting a prediction
  - 2. Trusting a model



## LIME – Desired Characteristics for Explainers

#### • Explanations must be interpretable

- Provide qualitative understanding between the input variables and the response
- Interpretability depends on the **target audience**
- To be meaningful an explanation must be locally faithful
  - Correspond to how the model behaves in the vicinity of the instance being predicted
- An Explainer should be model agnostic
- Provide global perspective
  - Establishing **trust** in the model

#### LIME – Goal

The overall goal of LIME is to identify an **interpretable** model over the *interpretable representation* that is **locally faithful** to the classifier.

## LIME – Interpretable Representation

#### Text classification:

- Feature:
  - Word Embedding
- Interpretable representation:
  - Binary vector indicating presence or absence of a word

Image classification:

- Feature:
  - Tensor with three color channels
- Interpretable representation:
  - Binary vector indicating the "presence" or "absence" of a super-pixel

#### LIME – Fidelity-Interpretablility Trade-off

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

#### LIME – Fidelity-Interpretablility Trade-off

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

- $g \in G$ , where G is a class of potentially **interpretable** models
  - Linear models
  - Decision trees

#### LIME – Fidelity-Interpretablility Trade-off

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

 $\Omega(g)$  is a measure of **complexity** 

- Linear model: number of non-zero weights
- Decision trees: depth of the tree
# LIME – Fidelity-Interpretablility Trade-off

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

 $f: \mathbb{R}^d \to \mathbb{R}$ , model to be explained

• In classification f(x) is the probability that x belongs to a certain class

# LIME – Fidelity-Interpretablility Trade-off

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

 $\pi_x(z)$ , **proximity** measure between an instance z to x

• To define locality around *x* 

# LIME – Fidelity-Interpretablility Trade-off

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

 $\mathcal{L}(f, g, \pi_x)$ , measure of **unfaithfulness** of g in approximating f in locality defined by  $\pi_x$ 

• Minimize  $\mathcal{L}(f, g, \pi_x)$  while keeping  $\Omega(g)$  low to ensure interpretability by humans

# LIME - Sampling for Local Exploration

- Black box model *f* (pink/blue)
- Explain single instance (bold red plus) with linear model
- Dashed line is learned locally faithful explanation
- NOT GLOBALLY FAITHFUL



# LIME - Superpixels



(a) Original Image (b) Explaining *Electric guitar* (c) Explaining *Acoustic guitar* (d) Exp [Riberio]

• Explanation of a single instance using superpixels

# **SP-LIME - Explaining Models**

Explanation of a single prediction is not sufficient to evaluate and assess trust in the model as a whole

- give global understanding of the model by explaining set of individual instances
- Use pick step
  - task of selecting B instances from X for the user to inspect
  - B is the budged of the user (time/patience)
  - pick a diverse, representative set of explanations to show the user
  - construct explanation matrix

# **SP-LIME – Explanation Matrix**

- Columns: features
- Rows: instances
- Cells contain local importance
- Global importance I
  - choose *I* such that features that explain many different instances have higher importance scores
- Pick problem consists of finding the set of Instances V that achieves the highest coverage



### §4 WHY

Why would one want to use visualization in deep learning?

### **§6 WHAT**

What data, features, and relationships in deep learning can be visualized?

## §8 WHEN

When in the deep learning process is visualization used?



### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

## §8 WHEN

When in the deep learning process is visualization used?



§5 WHO

Who would use and benefit from visualizing deep learning?

#### 7 HOW

How can we visualize deep learning data, features, and relationships?



Where has deep learning visualization been used?

### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

Computational Graph & Network Architecture Learned Model Parameters Individual Computational Units Neurons In High-dimensional Space Aggregated Information



When in the deep learning process is visualization used?





Who would use and benefit from visualizing deep learning?

#### 7 HOW

How can we visualize deep learning data, features, and relationships?



Where has deep learning visualization been used?

### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

Computational Graph & Network Architecture Learned Model Parameters Individual Computational Units Neurons In High-dimensional Space Aggregated Information

## §8 WHEN

When in the deep learning process is visualization used?

During Training After Training





Who would use and benefit from visualizing deep learning?

#### 7 HOW

How can we visualize deep learning data, features, and relationships?



Where has deep learning visualization been used?

### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

Computational Graph & Network Architecture Learned Model Parameters Individual Computational Units Neurons In High-dimensional Space Aggregated Information

## §8 WHEN

When in the deep learning process is visualization used?

During Training After Training





[Homann]

Who would use and benefit from visualizing deep learning?

Model Developers & Builders Model Users Non-experts



How can we visualize deep learning data, features, and relationships?



Where has deep learning visualization been used?

### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

Computational Graph & Network Architecture Learned Model Parameters Individual Computational Units Neurons In High-dimensional Space Aggregated Information

## §8 WHEN

When in the deep learning process is visualization used?

During Training After Training





# Who would use and benefit from visualizing deep learning?

Model Developers & Builders Model Users Non-experts



## How can we visualize deep learning data, features, and relationships?

Node-link Diagrams for Network Architecture Dimensionality Reduction & Scatter Plots Line Charts for Temporal Metrics



Where has deep learning visualization been used?

### §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability Debugging & Improving Models Comparing & Selecting Models Teaching Deep Learning Concepts

### §6 WHAT

What data, features, and relationships in deep learning can be visualized?

Computational Graph & Network Architecture Learned Model Parameters Individual Computational Units Neurons In High-dimensional Space Aggregated Information

## §8 WHEN

When in the deep learning process is visualization used?

During Training After Training





[Homann]

## Who would use and benefit from visualizing deep learning?

Model Developers & Builders Model Users Non-experts



# How can we visualize deep learning data, features, and relationships?

Node-link Diagrams for Network Architecture Dimensionality Reduction & Scatter Plots Line Charts for Temporal Metrics



## Where has deep learning visualization been used?

Application Domains & Models A Vibrant Research Community

# Future Development - Challenges

- Scalability
- Performance analysis
- Bias in representative examples
- Consensus on common visual approach

# Future Development - Opportunities

- Using expert knowledge
- Progressive visual analytics
- Advanced deep learning architectures
- Protection against adversarial attacks

Technical Term	Synonyms	Meaning
Neural Network	Artificial neural net, net	Biologically-inspired models that form the basis of deep learning; approximate functions dependent upon a large and unknown amount of inputs consisting of <i>layers</i> of <i>neurons</i>
Neuron	Computational unit, node	Building blocks of neural networks, entities that can apply activation functions
Weights	Edges	The trained and updated parameters in the neural network model that connect neurons to one another
Layer	Hidden layer	Stacked collection of <i>neurons</i> that attempt to extract features from data; a <i>layer's</i> input is connected to a previous <i>layer's</i> output
Computational Graph	Dataflow graph	Directed graph where nodes represent operations and edges represent data paths; when implementing <i>neural network</i> models, often times they are represented as these
Activation Functions	Transform function	Functions embedded into each <i>layer</i> of a <i>neural network</i> that enable the network represent complex non- linear decisions boundaries
Activations	Internal representation	Given a trained network one can pass in data and recover the <i>activations</i> at any <i>layer</i> of the network to obtain its current representation inside the network
Convolutional Neural Network	CNN, convnet	Type of <i>neural network</i> composed of convolutional <i>layers</i> that typically assume image data as input; these <i>layers</i> have depth unlike typical <i>layers</i> that only have width (number of <i>neurons</i> in a <i>layer</i> ); they make use of filters (feature & pattern detectors) to extract spatially invariant representations
Long Short-Term Memory	LSTM	Type of <i>neural network</i> , often used in text analysis, that addresses the vanishing gradient problem by using memory gates to propagate gradients through the network to learn long-range dependencies
Loss Function	Objective function, cost function, error	Also seen in general ML contexts, defines what success looks like when learning a representation, i.e., a measure of difference between a <i>neural network's</i> prediction and ground truth
Embedding	Encoding	Representation of input data (e.g., images, text, audio, time series) as vectors of numbers in a high- dimensional space; oftentimes reduced so data points (i.e., their vectors) can be more easily analyzed (e.g., compute similarity)
Recurrent Neural Network	RNN	Type of <i>neural network</i> where recurrent connections allow the persistence (or "memory") of previous inputs in the network's internal state which are used to influence the network output
Generative Adversarial Networks	GAN	Method to conduct unsupervised learning by pitting a generative network against a discriminative network; the first network mimics the probability distribution of a training dataset in order to fool the discriminative network into judging that the generated data instance belongs to the training set
Epoch	Data pass	A complete pass through a given dataset; by the end of one <i>epoch</i> , a <i>neural network</i> will have seen every datum within the dataset once

### Homann

# Sources

Adadi, A., & Berrada, M. (2018). Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6, 52138–52160. doi:10.1109/access.2018.2870052

Alicioglu, G., & Sun, B. (2022, February). A survey of visual analytics for Explainable Artificial Intelligence methods. *Computers & amp; Graphics, 102,* 502–520. doi:10.1016/j.cag.2021.09.002

Choo, J., & Liu, S. (2018). Visual Analytics for Explainable Deep Learning. Visual Analytics for Explainable Deep Learning. arXiv. doi:10.48550/ARXIV.1804.02527

Dragoni, M., Donadello, I., & Eccher, C. (2020, May). Explainable AI meets persuasiveness: Translating reasoning results into behavioral change advice. *Artificial Intelligence in Medicine*, *105*, 101840. doi:10.1016/j.artmed.2020.101840

Gunning, D., & Aha, D. W. (2019, June). DARPA's Explainable Artificial Intelligence Program. AI Magazine, 40, 44–58. doi:10.1609/aimag.v40i2.2850

Hohman, F., Kahng, M., Pienta, R., & Chau, D. H. (2019, August). Visual Analytics in Deep Learning: An Interrogative Survey for the Next Frontiers. *IEEE Transactions on Visualization and Computer Graphics*, 25, 2674–2693. doi:10.1109/tvcg.2018.2843369

La Rosa, B., Blasilli, G., Bourqui, R., Auber, D., Santucci, G., Capobianco, R., . . . Angelini, M. (2023, February). State of the Art of Visual Analytics for eXplainable Deep Learning. *Computer Graphics Forum*, *42*, 319–355. doi:10.1111/cgf.14733

Moradi, M., & Samwald, M. (2021, March). Post-hoc explanation of black-box classifiers using confident itemsets. *Expert Systems with Applications, 165,* 113941. doi:10.1016/j.eswa.2020.113941

Rai, A. (2019, December). Explainable AI: from black box to glass box. *Journal of the Academy of Marketing Science, 48*, 137–141. doi:10.1007/s11747-019-00710-5

Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. arXiv. doi:10.48550/ARXIV.1602.04938

Yuan, J., Chen, C., Yang, W., Liu, M., Xia, J., & Liu, S. (2021, March). A survey of visual analytics techniques for machine learning. *Computational Visual Media*, 7, 3–36. doi:10.1007/s41095-020-0191-7