

Visual Analytics for Explainable Machine Learning

Shixia Liu (刘世霞)

Tsinghua University (清华大学)

Recent ML Systems Achieve Superhuman Performance

AlphaGo beats Go human champ



Computer out-plays humans in "doom"



Deep Net outperforms humans in image classification

Autonomous search-and-rescue drones outperform humans



IBM's Watson destroys humans in jeopardy

	THINK		*		1	
-	\$300,000		\$1,000,000		\$200,000	
	Who is Stoker' fait as at it interp §1,000		Who Is Bram Sober7 \$ 17,973		who is INAMI STREEP? \$52.00	

DeepStack beats professional poker players



Deep Net beats human at recognizing traffic signs



Machine Learning as a Black Magic or Black Box



Most machine learning processes are NP-processes



3

Explainable AI – What Are We Trying To Do?





Explainable Artificial Intelligence - Darpa





Machine Learning Pipeline

Tasks

- Preprocess and clean the data



- Select and construct appropriate features
- Select an appropriate model family
- Optimize model hyperparameters
- Postprocess machine learning models
- Critically analyze the results obtained



Solution = Data + ML Expertise + Computation



Data is the King of Machine Learning!



AEVis, IEEE VIS 2018 (TVCG)

LabelInspect, IEEE VIS 2018 (TVCG)

Mengchen Liu*, Shixia Liu*, Hang Su[†], Kelei Cao*, Jun Zhu[†]

RS

-1911-





Adversarial Examples

 Intentionally designed to mislead a deep neural network (DNN) into making incorrect prediction



Technical Challenges

UNIVERSITY -1911--1911-

- Extract the <u>datapath</u> for adversarial examples
- Datapath visualization

A datapath

Hundreds of layers

Millions of neurons

Millions of connections



Datapath Extraction - Motivation

- Current methc Normal example
 - Most activate
- Problem
 - Misleading re highly recogr
- Reason
 - Neurons hav
 - Gap betweer



Adversarial example

Datapath Extraction - Formulation



- The critical neurons for a prediction: the neurons that highly contributed to the final prediction
- Subset selection
 - Keep the original prediction by selecting a minimized subset of neurons

$$N^{opt} = \underset{N_s \subseteq N}{\operatorname{arg\,min}} (p(x) - p(x; N_s))^2 + \lambda |N_s|^2$$

• Extend to a set of images X

$$N^{opt} = \underset{N_s \subseteq N}{\operatorname{arg\,min}} \sum_{x_k \in X} \left(p(x_k) - p(x_k; N_s) \right)^2 + \lambda |N_s|^2$$

Datapath Extraction - Solution



- Directly solving: time-consuming
 - NP-complete
 - Large search space due to the large number of neurons in a CNN

Quadratic approximation

Divide-and-conquer-based search space reduction

An accurate approximation in smaller search space

$$N^{opt} = \underset{N_s \subseteq N}{\operatorname{arg\,min}} \sum_{x_k \in X} \left(p(x_k) - p(x_k; N_s) \right)^2 + \lambda |N_s|^2$$

Datapath Extraction – Search Space Reduction



Original problem: 57.78 million dims



Datapath Extraction – Quadratic Approximation $F_{opt}^{i} = \underset{F_{s}^{i} \subseteq F^{i}}{\operatorname{arg\,min}} (p(x) - p(x; F_{s}^{i} \cup F^{-i}))^{2} + \lambda^{i} |F_{s}^{i}|^{2} \text{ Still NP}$ Reformulate $\mathbf{z}^i = [z_1^i, \dots, z_n^i]$ $\mathbf{z}_{opt}^{i} = \operatorname*{arg\,min}_{z_{j}^{i} \in \{0,1\}} (p(x) - p(x; \mathbf{z}^{i}))^{2} + \lambda^{i} (\sum_{j} z_{j}^{i})^{2}, \qquad z_{j}^{i} \text{ whether the } j\text{-th feature map in layer } i \text{ is critical map i$ $\mathbf{z}_{opt}^{i} = \underset{z_{j}^{i} \in [0,1]}{\operatorname{arg\,min}} (p(x) - p(x; \mathbf{z}^{i}))^{2} + \lambda^{i} (\sum_{j} z_{j}^{i})^{2}, \quad \text{Needs to calculate } \frac{\partial p}{\partial z_{j}^{i}} \text{ by BP in each iteration}$ 1. Bridge the gap between activation and prediction 2. Each element in Q \mathbf{a}_i : activation vector of approximately models the $\mathbf{z}_{opt}^{i} = \operatorname*{arg\,min}_{z_{j}^{i} \in [0,1]} \mathbf{z}^{i} \mathbf{Q}^{i} + \lambda^{i} \mathbf{I} (\mathbf{z}^{i})^{T} - 2q\mathbf{q}^{i} \cdot \mathbf{z}^{i},$ the j-th feature map interaction between $\mathbf{q}^{i} = [\mathbf{a}_{1}^{i} \cdot \frac{\partial p}{\partial \mathbf{a}_{1}^{i}} \Big|_{\mathbf{a}}, \cdots, \mathbf{a}_{n}^{i} \cdot \frac{\partial p}{\partial \mathbf{a}_{n}^{i}} \Big|_{\mathbf{a}}]$ feature map *j* and feature Quadratic optimization map k $Q = (\mathbf{q}^i)^T \mathbf{q}^i$

Datapath Visualization



SINGH,

-1911-

NERS/7





+noise

classified as a



Why?



Ф



•





.



A "diverging point" appears, where the activation similarity largely decreases.

 colobus, colobus monkey
 0 19

 guenon, guenon monkey
 0 44

4

-



guenon, guenon 0.44 monkey •



.

monkey



-



•





•













Want huge model capacity for large datasets





Analyzing the Training Processes of Deep Generative Models

Mengchen Liu, Jiaxin Shi, Kelei Cao, Jun Zhu, Shixia Liu Tsinghua University

IEEE VAST 2017 (IEEE TVCG)



Deep Generative Models (DGMs)





Training a DGM is Hard

- DGM often involves both deterministic functions and random variables
 - CNN: deterministic functions (e.g., convolution)
- DGM involves a top-down generative process and a bottom-up Bayesian inference process
 - CNN: a bottom-up process: input at the bottom layer
 -> high-level features -> outputs



Random variables

Challenges



- Handle a large amount of time series data
 - Millions of activations/gradients/weights in a DGM
- Identify the root cause of a failed training process
 - Errors may arise from multiple possible sources: abnormal training samples, inappropriate network structures, and lack of numerical stability in the library
 - Even when we can determine that the error is caused by the network structure, it is often difficult to locate the specific neurons

Our Solution



- A blue noise polyline sampling algorithm
 - Selects polyline samples the with blue-noise properties
 - Preserves outliers and reduce visual clutter

- A credit assignment algorithm
 - Discloses how other neurons contribute to the output of the neuron of interest





• Better understand and diagnose the training process of a DGM

Case Study: Debugging a Failed Training Process of a Variational Autoencoder (VAE)

- Autoencoder
 - Reconstruct their input with minimum information loss



- Variational autoencoder
 - Probabilistic version of an autoencoder
 - $-z_v$: a vector of random variables
 - za : a vector of real numbers

Dataset: CIFAR10 dataset Loss = NaN (10k-30k iterations) An example case: fails at 24,397



Data Flow: Output































Solution



• Trial 1:

- Replacing 🦛 with East, but the training failed again

- By the same analysis, we find another "bad" image
- in and in

• Trial 2:









Research Opportunities



- Human-in-the-loop visual analytics for practitioners
 - Existing deep learning models are data-driven
 - Combine human expert knowledge and deep learning techniques through interactive visualization
- Progressive visual analytics of deep learning models
 - The training of many deep learning models is time-consuming
 - Progressive visual analytics techniques are needed

Liu S, Wang X, Liu M, Zhu J. Towards better analysis of machine learning models: A visual analytics perspective. Visual Informatics. 2017 Mar 1;1(1):48-56. Jaegul Choo, Shixia Liu.Visual Analytics for Explainable Deep Learning. IEEE Computer Graphics and Applications, 2018.

Research Opportunities (Cont'd)



- Improving the robustness of deep learning models for secure artificial intelligence
 - Deep learning models are generally vulnerable to adversarial perturbations
 - Incorporate human knowledge to improve the robustness of deep learning models
- Reducing the size of the required training set
 - One-shot learning or zero-shot learning
- Visual analytics for advanced deep learning architectures
 - ResNet and DenseNet

References



- M. Liu, J. Shi, Z. Li, C. Li, J. Zhu, and S. Liu. Towards better analysis of deep convolutional neural networks. IEEE TVCG, 23(1):91–100, 2017
- M Liu, J Shi, K Cao, J Zhu, S Liu. Analyzing the Training Processes of Deep Generative Models. IEEE TVCG, , 24(1), 77-87, 2018.
- S. Liu, J. Xiao, J. Liu, X. Wang, J. Wu, & J. Zhu. Visual diagnosis of tree boosting methods. IEEE TVCG, 24(1), 163-173, 2018.
- Jaegul Choo, Shixia Liu. Visual Analytics for Explainable Deep Learning. IEEE Computer Graphics and Applications, 2018.
- S. Liu, X. Wang, M. Liu, and J. Zhu. Towards better analysis of machine learning models: A visual analytics perspective. Visual Informatics, 1(1):48– 56, 2017.

Acknowledgements

A UNIVERSITY

- Professors
 - Prof. Jun Zhu (Tsinghua), Prof. Hang Su (Tsinghua)
 - Prof. Jaegul Choo (Korea University)
- Students
 - Mengchen Liu, Kelei Cao, Changjian Chen, Fangxin Ouyang, Jiaxin Shi,
 Zhen Li, Chongxuan Li



Thanks a lot for your attention!

