Enhancing Simple Models by Interpretable Knowledge Transfer

Presenter:

1

Amit Dhurandhar, Research Staff Member, IBM T.J. Watson Research, Yorktown Heights, New York



Collaborators



Karthikeyan Shanmugam, Research Staff Member, IBM T.J. Watson Research, Yorktown Heights, New York



Ronny Luss, Research Staff Member, IBM T.J. Watson Research, Yorktown Heights, New York 3

Use Case: Fabrication engineers need to measure the amount of metal etched on each manufactured wafer (collection of chips).



Modeling: The engineers are comfortable using decision trees (e.g. CART) to model their problem because they understand the tool and know how to derive insights from the solution.

Challenge: Transfer information from complex models with higher accuracy to the decision trees in a manner that enhances performance and adds valuable insight to the engineers that is easily consumable by them.

Knowledge Transfer: General Idea





Domain Experts Preference: SMEs many times want to use a model they understand and hence trust. Our methods try to get the most out of these models by significantly enhancing their performance.

Small Data settings: Small client data where complex neural networks might overfit so simple models are preferred. However, possible to improve performance if we have a pretrained network on a large public or private corpora belonging to the same feature space.

Resource Constrained Settings: In memory and power constrained settings only small models can be deployed. Here improving small neural networks with the help of larger neural network can be useful.

Types of Knowledge Transfer (to Simple Models)

6

Fit to Black-Box Models Predictions

Knowledge Distillation

7

- Obtain soft predictions from a complex neural network (for different temperature scalings).
- Use them to regress a simpler neural network with cross-entropy loss.
- Model compression (Bucila et. al, 2006) where you fit to the logits or hard targets is a special case.



Distilling the Knowledge in a Neural Network. Hinton et. al. NeurIPS 2014 Interpreting Blackbox Models via Model Extraction. Bastani et. al. 2017

Sample Reweighting Methods

ProfWeight

- Find area under curve (AUC) based on confidence scores at each probe whose accuracy > α of the simple model for each example based on complex neural network.
- Weight each training example by corresponding AUC and retrain simple model such as a decision tree.



Improving Simple Models with Confidence Profiles. Dhurandhar et. al. NeurIPS 2018

Sample Reweighting Methods

SRatio

9

- Compute AUC like Profweight (also applicable to other classifiers than NNs).
- Divide by simple models confidence for that sample where the ratio is upper bounded by a tunable parameter (^β) which is its weight for retraining the simple model.



¹⁰ Local Explanations to Global Understanding



From local explanations to global understanding with explainable AI for trees. Lundberg et. al, Nat. Mach. Intl. 2020

Model Agnostic Multilevel Explanations. Ramamurthy et. al, NeurIPS 2020 Learning Global Transparent Models Consistent with Local Contrastive Explanations. Pedapati et. al, NeurIPS 2020 11

Sample Reweighting Methods

Sample Reweighting Methods

What are Probes?

12

- Probes are logistic classifiers i.e. a linear model with bias followed by softmax.
- We train them on intermediate representations of a neural network (NN) keeping the weights of the NN fixed.



Sample Reweighting Methods

ProfWeight

13

- Find area under curve (AUC) based on confidence scores at each probe whose accuracy > α of the simple model for each example based on complex neural network.
- Weight each training example by corresponding AUC and retrain simple model such as a decision tree.



Improving Simple Models with Confidence Profiles. Dhurandhar et. al. NeurIPS 2018 https://www.kdnuggets.com/2020/09/deep-learning-dream-accuracy-interpretability-single-model.html

Algorithm 1 ProfWeight

Input: k unit neural network \mathcal{N} , learning algorithm for simple model $\mathcal{L}_{\mathcal{S}}$, dataset D_N used to train \mathcal{N} , dataset $D_S = \{(x_1, y_1), ..., (x_m, y_m)\}$ to train a simple model and margin parameter α .

1) Attach probes $P_1, ..., P_k$ to the k units of \mathcal{N} .

2) Train probes based on D_N and obtain errors $e_1, ..., e_k$ on D_S . {There is no backpropagation of gradients here to the hidden units/layers of \mathcal{N} .}

3) Train the simple model $\mathcal{S} \leftarrow \mathcal{L}_{\mathcal{S}}(D_S, \beta, \vec{1}_m)$ and obtain its error e_S . { \mathcal{S} is obtained by unweighted training. $\vec{1}_m$ denotes a *m* dimensional vector of 1s.}

4) Let $I \leftarrow \{u \mid e_u \leq e_S - \alpha\} \{I \text{ contains indices of all probes that are more accurate than the simple model <math>S$ by a margin α on D_S .

5) Compute weights w using Algorithm 2 or 3 for AUC or neural network, respectively.

6) Obtain simple model $S_{\mathbf{w},\beta} \leftarrow \mathcal{L}_{S,\beta}(D_S,\beta,w)$ {Train the simple model on D_S along with the weights w associated with each input.}

return $S_{w,\beta}$

¹⁵ Ideally (ProfWeight Version 2),



- 1. Minimize loss of simple model (βs) for current weights.
- 2. For given parameters of the simple model minimize over the weights (w).
- 3. Repeat steps 1 and 2 until convergence or max iterations.



Problem: Weights may all go to zero

Solution: Add regularization $\mathcal{R}(\mathbf{w}) = (\frac{1}{m} \sum_{i=1}^{m} w_i - 1)^2$

$$S^* = S_{\mathbf{w}^*,\beta^*} = \min_{\mathbf{w}\in\mathcal{C}} \min_{\beta\in\mathcal{B}} E\left[\lambda\left(S_{\mathbf{w},\beta}(x),y\right)\right]$$

Subject to:
$$\mathcal{R}(\mathbf{w}) = (\frac{1}{m} \sum_{i=1}^{m} w_i - 1)^2$$

¹⁷ Theoretical Justification

Theorem 1. $\mathbb{E}_{P_M(\mathbf{x}|y)}[w_{M'}(\mathbf{x},y)] = 1$ implies that $w_{M'}(\mathbf{x},y) = \frac{P(\mathbf{x}|y)}{P_M(\mathbf{x}|y)}$.

It is well-known that the error of the performance of the best classifier (Bayes optimal) on a distribution of class mixtures is the total variance distance between them. That is:

Lemma 2. [8] The error of the Bayes optimal classifier trained on a uniform mixture of two class distributions is given by: $\min_{\theta} \sum \mathbb{D}[L_{\theta}(x, y)] = \frac{1}{2} - \frac{1}{2}D_{\text{TV}}(P(\mathbf{x}|y=1), P(\mathbf{x}|y=0))$ where $L(\cdot)$ is the 0, 1 loss function and θ is parameterization over a class of classifiers that includes the Bayes optimal classifier. D_{TV} is the total variation distance between two distributions. $P(\mathbf{x}|y)$ are the class distributions in dataset D.

From Lemma 2 and Theorem 1, it is clear that:

$$\min_{M',\theta \text{ s.t. } \mathbb{E}_{\bar{D}_{\text{train}}}[w'_{M}]=1} \mathbb{E}_{\bar{D}}[w_{M'}(\mathbf{x}, \mathbf{y})L_{\theta}(\mathbf{x}, y)] = \frac{1}{2} - \frac{1}{2}D_{\text{TV}}(P(\mathbf{x}|y=1), P(\mathbf{x}|y=0))$$
(2)

¹⁸ Experiments

Dataset: CIFAR-10

Complex Model: 18 Unit Resnet

Simple Models:

3 Unit Resnet5 Unit Resnet7 Unit Resnet9 Unit Resnet

Performance Improvement: 3-4%

	SM-3	SM-5	SM-7	SM-9
Standard	$73.15(\pm 0.7)$	$75.78(\pm 0.5)$	$78.76(\pm 0.35)$	$79.9(\pm 0.34)$
ConfWeight	$76.27 (\pm 0.48)$	$78.54 (\pm 0.36)$	81.46(±0.50)	$82.09 (\pm 0.08)$
Distillation	$65.84(\pm 0.60)$	$70.09 (\pm 0.19)$	$73.4(\pm 0.64)$	$77.30 (\pm 0.16)$
ProfWeight ^{ReLU}	77.52 (±0.01)	$78.24(\pm 0.01)$	$80.16(\pm 0.01)$	$81.65 (\pm 0.01)$
ProfWeight ^{AUC}	$76.56 (\pm 0.62)$	$79.25(\pm 0.36)$	$81.34(\pm 0.49)$	82.42 (±0.36)



¹⁹ Experiments



20

Conceptual Contribution: Upweighting simpler examples (w.r.t. the complex model) helps (somewhat opposite intuition to boosting, but not really ...).

Algorithmic Contribution: Proposed a method, called ProfWeight, based on the above conceptual idea that actually accomplishes improvement using Probes.

Theoretical Contribution: We show that the alternating optimization scheme has justification (although hard to control in practice).

Sample Reweighting Methods

SRatio

21

- Compute AUC like Profweight (also applicable to other classifiers than NNs).
- Divide by simple models confidence for that sample where the ratio is upper bounded by a tunable parameter (^β) which is its weight for retraining the simple model.



Enhancing Simple Models by Exploiting What They Already Know. Dhurandhar et. al. ICML 2020

Intuition for SRatio

Covariate Shift

<u>Classifiers</u>

p is source, q is target

p(y|x) = q(y|x)

But $p(x) \neq q(x)$

Typical Soln: Weight by q(x)/p(x) p is simple, q is complex

p(x) = q(x)

But $p(y|x) \neq q(y|x)$

Our (analogous) Soln: Weight by q(y|x)/p(y|x)

(Inverse covariate shift?)

²³ Theory

Lemma 3.1. Let $p_{\theta}(y|x)$ be the softmax scores on a specific model θ from simple model space Θ . Let $\theta^* \in \Theta$ be the set of simple model parameters that is obtained from a given learning algorithm for the simple model on a training dataset. Let $p_c(y|x)$ be a pre-trained complex classifier whose loss is smaller than θ^* on the training distribution. Let $\beta \ge 1$ be a scalar clip level for the ratio $p_c(y|x)/p_{\theta^*}(y|x)$. Then we have:

$$\mathbb{E}[-\log p_{\theta}(y|x)] \leq \mathbb{E}\left[\max\left(1, \min\left(\frac{p_{c}(y|x)}{p_{\theta^{*}}(y|x)}, \beta\right)\right) \log\left(\frac{1}{p_{\theta}(y|x)}\right)\right] - \mathbb{E}\left[\log\left(\min\left(\frac{p_{c}(y|x)}{p_{\theta^{*}}(y|x)}, \beta\right)\right)\right] + \log(\beta).$$
(1)

²⁴ Probes: To judge hardness of examples in NNs



Improving Simple Models with Confidence Profiles. Dhurandhar et. al. NeurIPS 2018

Generalization of Probes: Graded Classifiers

Definition (δ -graded) Let $X \times Y$ denote the input-output space and p(x, y) the joint distribution over this space. Let $\zeta_1, \zeta_2, ..., \zeta_n$ denote classifiers that output the prediction probabilities for a given input $x \in X$ for the most probable (or true) class $y \in Y$ determined by p(y|x). We then say that classifiers $\zeta_1, \zeta_2, ..., \zeta_n$ are δ -graded for some $\delta \in (0, 1]$ and a (measurable) set $Z \subseteq X$ if $\forall x \in Z$, $\zeta_1(x) \leq \zeta_2(x) \leq \cdots \leq \zeta_n(x)$, where $\int_{x \in Z} p(x) \geq \delta$.

Deep Neural Networks: (Linear) probes to intermediate representations.

Boosted Trees: Average predictions of first k trees, first 2k trees, ...

Random Forests: Order trees by accuracy and average predictions of first k trees, first 2k trees, ...

Other Models: Taking increasing order Taylor approximations or do functional decomposition.

25

²⁶ SRatio Algorithm

Algorithm 1 Our proposed method SRatio.

Input: n (graded) classifiers $\zeta_1, ..., \zeta_n$, learning algorithm for simple model \mathcal{L}_S , dataset D_S of cardinality N, performance gap parameter γ and maximum allowed ratio parameter β .

Train simple model on D_S, S ← L_S(D_S, 1_N) and compute its (average) prediction error ε_S.{Obtain initial simple model where each input is given a unit weight.}
Compute (average) prediction errors ε₁, ..., ε_n for the n graded classifiers and store the ones that are at least γ more accurate than the simple model i.e. I ← {i ∈ {1, ..., n} | ε_S - ε_i ≥ γ}
Compute weights for all inputs x as follows: w(x) = ∑_{i∈I} ζ_{i(x)}/mS(x), where m is the cardinality of set I and S(x) is the prediction probability/score for the true class of the simple model.
Set w(x) ← 0, if w(x) > β {Limit the importance of extremely hard examples for the simple model.}
Retrain the simple model on the dataset D_S with the corresponding learned weights w, S_w ← L_S(D_S, w)

6) Return S_w

Comparison with other transfer methods

27

Method	Complex Model	Simple Model	Models Leveraged
Distillation	NN	NN	Complex
ProfWeight	NN	Any	Complex
SRatio	Any	Any	Complex and Simple

²⁸ Results

Complex Models

- Boosted Trees
- Random Forests
- ResNet 18

Simple Models

- Decision tree
- SVM
- ResNet 3, 5, 7

Table 1. Dataset characteristics, where N denotes dataset size and d is the dimensionality.

Dataset	N	d	# of Classes
Ionosphere	351	34	2
Ovarian Cancer	216	4000	2
Heart Disease	303	13	2
Waveform	5000	40	3
Human Activity	10299	561	6
Musk	6598	166	2
CIFAR-10	60000	32×32	10

²⁹ Results

	Complex	CM	Simple	SM Distill-proxy 1		ConfWeight	SRatio
Dataset	Model	Error	Model	Error	Error (SM)	Error (SM)	Error (SM)
			Tree	10.95	10.95	11.42	8.57^{*}
	Boosted	8.10		±0.4	±0.4	±0.8	± 0.5
	Trees	±0.4	SVM	12.38	11.90	11.90	10.47
Ionoenhora				±0.6	±0.6	± 0.6	± 0.5
ionosphere			Tree	10.95	10.95	11.42	10.42
	Random	6.19		±0.4	±0.4	± 0.4	± 0.1
	Forest	±0.4	SVM	12.38	12.38	12.38	11.42
				±0.6	±0.6	± 0.6	± 0.3
			Tree	15.62	15.62	15.62	15.62
	Boosted	4.68		± 0.8	±0.8	± 1.0	± 0.5
	Trees	±0.4	SVM	1.56	1.56	1.56	1.56
Overien Concer	1			±0.4	± 0.4	± 0.4	± 0.4
Ovarian Cancer			Tree	15.62	15.62	14.06	14.04
	Random	6.25		±0.8	±0.8	± 0.1	± 0.1
	Forest	±0.8	SVM	1.56	1.56	1.56	1.56
				± 0.4	± 0.4	± 0.4	± 0.4
			Tree	23.88	22.77	23.33	22.77
	Boosted	15.55		±0.7	± 0.1	± 0.3	± 0.2
	Trees	±0.6	SVM	17.22	16.67	17.22	16.77
Hoart Diseases				±0.2	± 0.3	± 0.2	± 0.2
ficant Disease			Tree	23.88	23.88	25.55	22.77
	Random	15.88		± 0.7	± 0.7	± 0.5	± 0.3
	Forest	±0.6	SVM	17.22	17.22	16.67	16.67
				±0.2	±0.2	± 0.3	± 0.2

³⁰ Results

			Tree	25.43	25.06	25.10	25.06
	Boosted	12.96		± 0.2	± 0.1	± 0.1	± 0.1
	Trees	±0.1	SVM	14.70	15.33	14.70	13.72
Wayafarm				± 0.2	±0.0	± 0.2	± 0.2
waveronn			Tree	25.43	25.43	25.43	25.06
	Random	10.90		± 0.2	± 0.2	± 0.2	± 0.1
	Forest	± 0.1	SVM	14.70	14.33	14.30	12.72
				± 0.2	±0.0	± 0.2	± 0.5
			Tree	7.93	7.93	7.86	7.15
	Boosted	6.32		± 0.2	± 0.1	± 0.2	± 0.1
	Trees	±0.0	SVM	14.56	15.85	13.92	13.92
Human Activity				± 0.1	± 0.1	± 0.1	± 0.2
Recognition			Tree	7.93	7.23	7.21	6.67
	Random	2.34		± 0.2	± 0.1	± 0.1	± 0.0
	Forest	± 0.0	SVM	14.56	13.92	14.24	13.92
				± 0.1	± 0.1	± 0.1	± 0.1
			Tree	4.49	6.11	4.45	4.06*
	Boosted	4.06		± 0.1	± 0.1	± 0.1	± 0.1
	Trees	± 0.1	SVM	6.11	6.29	6.41	5.48
Muck				± 0.1	± 0.1	± 0.1	± 0.1
WIUSK			Tree	4.49	4.49	4.47	3.89
	Random	2.45		± 0.1	± 0.1	± 0.1	± 0.1
	Forest	± 0.1	SVM	6.11	6.16	5.96	5.53
				± 0.1	± 0.1	± 0.1	± 0.1

³¹ Results



³² Results

³³ Results

Table 3. Below we see the % of training points whose weights based on SRatio have changed by > 1% compared to just considering the complex model, i.e., ignoring simple model confidences in Step 3 of algorithm 1. This depicts the impact of considering the simple models confidences in the weighting, which is our main conceptual contribution. Results are averaged over both complex models.

Ionos	phere	00	C	H	D	Wave	eform	H	AR	M	usk
Tree	SVM	Tree	SVM	Tree	SVM	Tree	SVM	Tree	SVM	Tree	SVM
37.40	90.65	8.55	6.57	92.02	99.06	44.9	99.7	5.9	29.45	7.5	13.93
(± 2.4)	(± 1.1)	(± 0.2)	(±1)	(± 1.8)	(± 0.1)	(± 2.4)	(± 0.1)	(± 0.9)	(± 1.3)	(± 0.5)	(± 0.7)

³⁴ Results on CIFAR-10

Table 4. Below we observe the averaged accuracies (%) of simple models SM-3 (3 Res units), SM-5 (5 Res units) and SM-7 (7 Res units) trained with various weighting methods and distillation. The complex model achieved 84.5% accuracy. Statistically significant best results are indicated in bold.

	SM-3	SM-5	SM-7
Standard	73.15 (± 0.7)	75.78 (±0.5)	78.76 (±0.35)
ConfWeight	76.27 (±0.48)	78.54 (±0.36)	81.46 (±0.50)
Distillation	65.84 (±0.60)	70.09 (±0.19)	73.4 (±0.64)
ProfWeight	76.56 (±0.51)	79.25 (±0.36)	81.34 (±0.49)
SRatio	77.23 (±0.14)	80.14 (±0.22)	81.89 (±0.28)

Takeaways SRatio

Conceptual Contribution: Simple models can be useful to improve themselves (without increasing complexity).

Algorithmic Contribution: Proposed a method, called SRatio, based on the above conceptual idea that actually accomplishes that.

Theoretical Contribution: We show that our weighting scheme is principled as it is a (tight) upper bound on the simple models expected loss.

Formalization: We propose a notion of delta-graded classifiers which is a formalization and generalization of the idea of probes.

1. Experiments on more domains and datasets.

2. Try to more deeply understand why the ratio weighting works.

3. Develop other adaptive methods that perform better (viz. multihop methods (blog to come out Oct. end in kdnuggets)).

4. Develop methods that are also resilient to hardware faults.

Even if you desire just a simple model...

build the most accurate model you can and transfer information to the simple model

Thank you

Other Directions: Contrastive Explanations

Anomaly Attribution with Likelihood Compensation. Ide et. al. 2020

Generating Contrastive Explanations with Monotonic Attribute Functions. Russ et. al. arxiv 2019

Other Directions

40

- **1. Meta-Methods Local to Global:** a) From LIME-like techniques can build a tree of explanations and b) From contrastive/counterfactual explanations can generate global Boolean features which can be used to train simple models.
- 2. Modeling and Handling Cognitive Biases: We have done some recent work on modeling cognitive biases and user tested an adaptive time based strategy to debias.
- 3. Explaining Similarity Based Classifiers: Feature based and analogy based (new type) explanations.
- 4. Causal Disentanglement: Uncovering latents with particular causal structures.
- 5. Generating Neighborhoods with Business Rules: For local explainability techniques how do we generate good neighborhoods that represent business rules well.
- 6. Interpretable Strategies: Given a complicated strategy by an agent can we disentangle it into understandable, but accurate sub-goals.
- 7. Explaining Representations: Given an arbitrary representation can I explain it in terms of an understandable representation.