

Harnessing Deep Neural Networks with Logic Rules

Reporter: HE WEINAN

School of Data and Computer Science
Sun Yat-sen University

cura4ho@gmail.com

May 8, 2017

- Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard Hovy, and Eric Xing. 2016. *Harnessing deep neural networks with logic rules*. In *Proc. of ACL*.
- http://www.cs.cmu.edu/~zhitingh/data/acl16harnessing_slides.pdf

Cited by

- Alashkar et al. (2017). Examples-Rules Guided Deep Neural Network for Makeup Recommendation. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Hu et al. (2016). Deep neural networks with massive learned knowledge. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

- 1 Motivation
- 2 Neural Network
- 3 Knowledge Distillation
 - Soft Logic
 - Distillation
- 4 Applications and Experiments
 - Experiments
- 5 Conclusion

Deep neural networks (DNN) is a powerful mechanism for learning patterns from massive data, achieving great performance on many problem domains such as *image classification* and *machine translation*.

But they still have limitations:

- Relying heavily on massive labeled training data
- The resulting parameters are often uninterpretable
- Hard to encode human knowledge and intention

On contrary, humans learn from

- concrete examples just like DNN do
- general **knowledge**

Logic rules is a flexible way to express structured knowledge. Integrating logic rules into DNNs might help to transfer human intention and domain knowledge the DNN models.

The authors propose an **knowledge distillation** framework to train a neural network using both examples and rules. A "teacher" is iteratively constructed to "teach" the student network.

Firstly, we need to know how traditional *neural networks* works. Consider an example of a classification problem.

Classification Problem: Loan Defaults

Given a *training set* of records of clients who take a loan, we'd like to predict whether an unseen client will default (not paying the loan).

Tid	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Label
Class

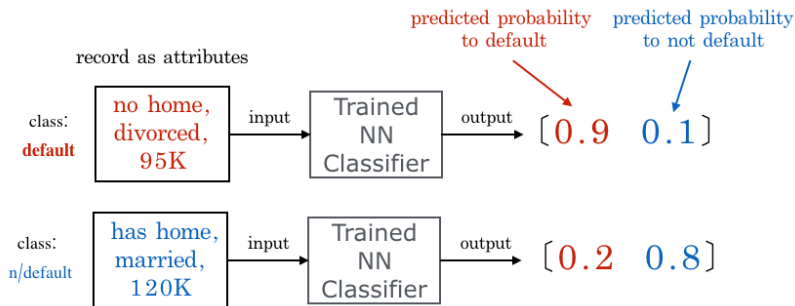
not default
default

Training set for predicting borrowers who will default on loan payments.

If the client *John* has no home, is divorced and makes 50K a year, how likely will he default?

A **neural network** (NN) classifier is a computational model that classifies client records into two **classes**: default or not default. These are called class **labels**

For example, suppose we have already trained such a NN classifier, then for the two training example



For an input record x with its actual label y , a NN classifier is a function that maps x into a soft prediction vector σ_θ :

$$f_\theta(x) = \sigma_\theta(x)$$

$\theta = (\theta_1, \theta_2, \dots, \theta_n)$: the **parameters** of the NN

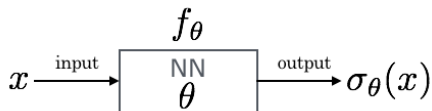
$x = (a_1, a_2, a_3)$: **input variable**, a_i is the i^{th} attribute of x

$\sigma_\theta = (p(y = \text{default}|x), p(y = \text{not default}|x))$: **soft prediction vector**

Note that the prediction vector $\sigma_\theta(x)$ forms a probability distribution $p_\theta(y|x)$.

Parameters or Weights

The **parameters** θ , or **weights** of the NN are the result of training on labeled data. θ determines the performance of the NN.



Example (Classifying A Record)

For the record $x = (\text{no home, divorced, 95K})$ which actually defaults,

$$f_{\theta}(x) = \sigma_{\theta}(x) = (0.9, 0.1)$$

The NN predicts that the probability of x default is 0.9. If the NN is perfect, the prediction should be $(1, 0)$.

We use a *label vector* $(1, 0)$ to represent default, $(0, 1)$ for not default.

To measure the performance: compare the prediction $(0.9, 0.1)$ with the actual label $y = (1, 0)$.

Loss Function

The **loss function** $\ell(y, \sigma_{\theta}(x))$ measures the difference between the prediction vector and the label vector for every training example.

To improve the performance: change parameters θ to minimize the loss function. (*Gradient Descent*)

Objective

Iteratively change parameters θ to minimize the loss function for all training examples.

Goal

Find θ that minimize the loss ℓ on N training examples.

$$\theta = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{y}_n, \sigma_{\theta}(\mathbf{x}_n))$$

Input: labeled training data $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$

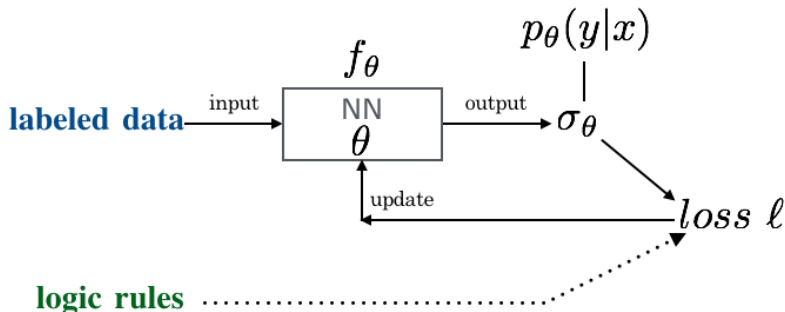
Output: parameters (weights) θ and prediction p_{θ}

Initialize the params $\theta^{(0)}$, at iteration t , **repeat:**

- ① Sample a subset $(\mathbf{X}, \mathbf{Y}) \subset \mathcal{D}$
- ② Compute the loss on (\mathbf{X}, \mathbf{Y})
- ③ Update the params $\theta^{(t)}$ from the previous params $\theta^{(t-1)}$

until $\theta^{(t)}$ converges, then let $\theta = \theta^{(t)}$.

Note: the traditional NNs ONLY train θ using labeled training data.



Idea

- 1 Encode the knowledge rules using **Soft Logic**
- 2 Change the loss function to incorporate logic rules
- 3 When updating, θ is influenced by the rules

- In *first order logic*, propositions are evaluated to truth value *True* or *False*.
- In *soft logic*, propositions and groundings (expressions with all variables being instantiated) evaluate to a continuous truth value from the interval $[0, 1]$.

The Boolean operators are extended as:

$$\begin{aligned}A \& B &= \max\{A + B - 1, 0\} \\ A \vee B &= \min\{A + B, 1\} \\ \neg A &= 1 - A\end{aligned}\tag{1}$$

Denote

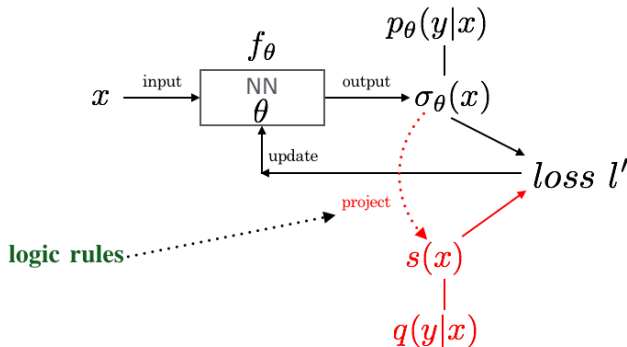
- R_l as the l^{th} rule over the input - target label space,
- λ_l as confidence level of rule R_l , with $\lambda_l = \infty$ indicating all groundings must be true.

Example

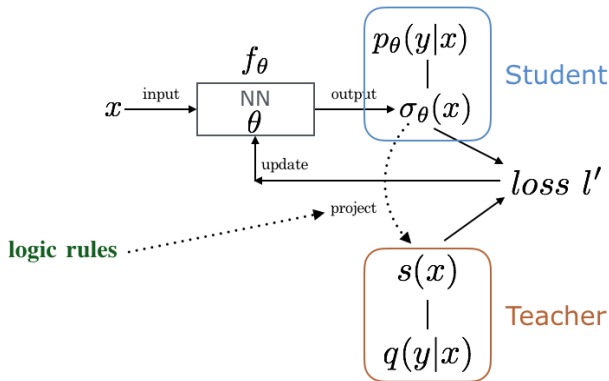
Rule: "People who earn less than 100K and are not married will default"

$$R \equiv (less_{100}(\mathbf{x}) \ \& \ \neg married(\mathbf{x})) \rightarrow default(\mathbf{x})$$

Now consider using the extracted rules to influence the loss function.



- NN maps the input x into prediction vector $\sigma_\theta(x)$ that forms probability $p_\theta(y|x)$.
- Project $\sigma_\theta(x)$ into $s(x)$, a rule-regularized prediction that forms probability $q(y|x)$.
- Then $q(y|x)$ contains information of the rules.



Example (Rule-regularized Projection)

For a client x that $less_{100}(x) \ \& \ \neg married(x)$,

the NN outputs a prediction $\sigma_\theta(x) = [0.7, 0.3]$,

the projected prediction might be $s(x) = [0.8, 0.2]$.

The new loss function l' will also consider the difference between $\sigma_\theta(x)$ and $s(x)$.

- Prediction $p_{\theta}(y|\mathbf{x})$ based only on examples (*student*)
- Prediction $q(y|\mathbf{x})$ also considers rules (*teacher*)
- **Distillation**: training θ to imitate the outputs that consider the rules. (train the *student* p_{θ} to imitate the *teacher* q)

at iteration t :

$$\theta^{(t+1)} = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N (1 - \pi) \ell(\mathbf{y}_n, \sigma_{\theta}(\mathbf{x}_n)) + \pi \ell(\mathbf{s}_n^{(t)}, \sigma_{\theta}(\mathbf{x}_n)),$$

Annotations:

- green text: true hard label (points to \mathbf{y}_n)
- blue text: soft prediction of student p_{θ} (points to $\sigma_{\theta}(\mathbf{x}_n)$ in the first term)
- blue text: balancing parameter (points to π)
- orange text: soft prediction of the teacher network q (points to $\mathbf{s}_n^{(t)}$)

With the set of rules $\mathcal{R} = \{(R_l, \lambda_l)\}_{l=1}^L$, the goal is to construct the teacher q at each iteration from p_θ such that:

- 1 q fits the rules
- 2 q stays close to p_θ

Essentially we can compute such a q by solving the optimization problem:

$$\begin{aligned} \min_{q, \xi \geq 0} & \quad \boxed{\text{KL}(q \| p_\theta(\mathbf{Y} | \mathbf{X}))} + C \sum_l \xi_l \quad \text{close to } p_\theta \\ \text{s.t.} & \quad \boxed{\lambda_l (1 - \mathbb{E}_q[r_l(\mathbf{X}, \mathbf{Y})]) \leq \xi_l} \\ & \quad l = 1, \dots, L \quad \text{rule constraints} \end{aligned}$$

The closed-form solution:

$$q^*(\mathbf{Y} | \mathbf{X}) \propto p_\theta(\mathbf{Y} | \mathbf{X}) \exp \left\{ - \sum_l C \lambda_l (1 - r_l(\mathbf{X}, \mathbf{Y})) \right\}$$

Goal

Given \mathcal{R} , find θ that minimize the loss ℓ' on N training examples.

$$\theta = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N (1 - \pi) \ell(\mathbf{y}_n, \sigma_{\theta}(\mathbf{x}_n)) + \pi \ell(\mathbf{s}_n(\mathbf{x}_n), \sigma_{\theta}(\mathbf{x}_n))$$

Input: Training data $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$,
The rule set $\mathcal{R} = \{(R_l, \lambda_l)\}_{l=1}^L$,
 π : imitation parameter, C : regularization strength

Output: student network p_{θ} and teacher network q

Initialize the parameters θ , then **repeat**:

- ① Sample a subset $(\mathbf{X}, \mathbf{Y}) \subset \mathcal{D}$
- ② Compute student p_{θ}
- ③ Construct teacher network q
- ④ Transfer knowledge into p_{θ} by updating the params θ

until θ converges

- Sentiment classification: sentence \rightarrow positive / negative
- Base neural network: CNN with accuracy about 87% (SST2)

For example, the sentence "This soup is good, but I don't like it" has negative sentiment although the first part seems to be positive.

Notice "but" changes the sentiment, we thus consider a simple rule:

If a sentence S has a structure " A -but- B ", then

$$\textit{sentiment}(S) = \textit{sentiment}(B)$$

Using the above rule and the knowledge distillation method, the authors observe a boost on accuracy.

Accuracy (%) of Sentiment Classification with all labeled data

	Model	SST2	MR	CR
1	CNN (Kim, 2014)	87.2	81.3 \pm 0.1	84.3 \pm 0.2
2	CNN-Rule- <i>p</i>	88.8	81.6 \pm 0.1	85.0 \pm 0.3
3	CNN-Rule- <i>q</i>	89.3	81.7\pm0.1	85.3\pm0.3
4	MGNC-CNN (Zhang et al., 2016)	88.4	–	–
5	MVCNN (Yin and Schutze, 2015)	89.4	–	–
6	CNN-multichannel (Kim, 2014)	88.1	81.1	85.0
7	Paragraph-Vec (Le and Mikolov, 2014)	87.8	–	–
8	CRF-PR (Yang and Cardie, 2014)	–	–	82.7
9	RNTN (Socher et al., 2013)	85.4	–	–
10	G-Dropout (Wang and Manning, 2013)	–	79.0	82.1

Row 2 and 3 are networks enhanced by the knowledge distillation, they outperform the base CNN and achieve better or comparable result with others.

Accuracy (%) of Sentiment Classification on SST2 dataset with varying sizes size of labeled data and semi-supervised learning.

	Data size	5%	10%	30%	100%
1	CNN	79.9	81.6	83.6	87.2
2	-Rule- p	81.5	83.2	84.5	88.8
3	-Rule- q	82.5	83.9	85.6	89.3
4	-semi-PR	81.5	83.1	84.6	—
5	-semi-Rule- p	81.7	83.3	84.7	—
6	-semi-Rule- q	82.7	84.2	85.7	—

Row 1 to 3 use only labeled data, while row 4 to 6 use the remaining data as unlabeled examples to train the NN. Row 2, 3, 5, 6 are results of the authors.

Contributions:

- Combine DNNs with logic rules to integrate knowledge
- Iteratively transfer knowledge
- General applicability: CNNs/RNNs
- Better performance using only one rule

Drawbacks:

- We have to feed rules to the NN
- For many tasks other than NLP (such as CV), clear rules are difficult to come up with
- Can more rules lead to better performance?

Possible future work:

- Derive rules from data?
- In another domain of problems? (QA / Image QA / Vision)