

# Deciphering Raw Data in Neuro-Symbolic Learning with Provable Guarantees

Lue Tao, Yu-Xuan Huang, Wang-Zhou Dai, Yuan Jiang

Shi-Jie Sang  
April 18, 2025

# Conventional Supervised Learning

Let  $\mathcal{X} \subset \mathbb{R}^d$  be the input space and  $\mathcal{Y} = \{0, 1, \dots, c - 1\}$  be the label space. The objective is to learn a mapping  $h : \mathcal{X} \rightarrow \mathbb{R}^c$  that minimises the expected risk:

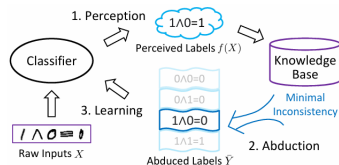
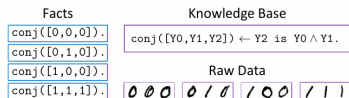
$$\mathcal{R}(h) = \mathbb{E}_{p(x,y)} \ell(h(x), y),$$

where  $\ell : \mathbb{R}^c \times (y) \rightarrow \mathbb{R}$  is a loss function that measures how well the classifier perceives an input.

# Neuro-Symbolic Learning

The raw inputs  $X = [x_0, x_1, \dots, x_{m-1}]$  are given, while their ground-truth  $Y = [y_0, y_1, \dots, y_{m-1}]$  are not observable. Instead, we only know that the logical facts grounded by the labels are compatible with a given knowledge base.

# Example

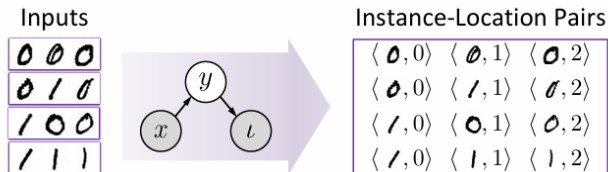


$$\mathcal{R}_{NeSy}(h) = \mathbb{E}_{p(X,\tau)} \mathcal{L}(X, \bar{Y}; h), s.t. B \cup \bar{Y} \models \tau$$

# Heuristics to search for the most likely labels

- Maximal probability:  $\bar{Y} = \operatorname{argmax}_Y \hat{p}(Y | X)$
- Minimal distance:  $\bar{Y} = \operatorname{argmin}_Y ||Y - \hat{Y}||$
- Random: Randomly select an element  $Y$  from the candidate set
- .....

# Instance-Location Pairs



From  $n$  sequences of unlabelled data  $\{X^{(i)}\}_{i=0}^{n-1}$ , a total of  $mn$  instance-location pairs  $\{\langle x^{(i)}, \iota^{(i)} \rangle\}_{i=0}^{mn-1}$  will be obtained.

# Location-Based Risk

**Assumption.**  $p(\iota \mid y, x) = p(\iota \mid y)$ .

$$p(\iota = k \mid x) = \sum_{j=0}^{c-1} Q_{jk} \cdot p(y = j \mid x),$$

where  $Q_{jk} = p(\iota = k \mid y = j)$ .

In practice, we let  $q(x) = Q^T g(x)$  and interpret  $g(x)$  as probabilities via  $g_j(x) = \exp(h_j(x)) / \sum_{i=0}^{c-1} \exp(h_i(x))$ . Then we minimise the Location-based risk:

$$\mathcal{R}_L(h) = \mathbb{E}_{p(x, \iota)} \ell(q(x), \iota).$$

# Example

Facts	Knowledge Base
conj0([0,0]).	conj([Y0,Y1,Y2]) ← See Figure 2.
conj0([0,1]).	conj0([Y0,Y1]) ← conj([Y0,Y1,0]).
conj0([1,0]).	conj1([Y0,Y1]) ← conj([Y0,Y1,1]).
conj1([1,1]).	

A set of  $mn$  triplets  $\{\langle x^{(i)}, \iota^{(i)}, \tau^{(i)} \rangle\}_{i=0}^{mn-1}$  will be obtained. For conciseness, we use  $\tilde{y}$  to denote  $\langle \iota, \tau \rangle$ . Then, the triplets above becomes  $\{\langle x^{(i)}, \tilde{y}^{(i)} \rangle\}_{i=0}^{mn-1}$ . Similar to before, we have

$$p(\tilde{y} = o \mid x) = \sum_{j=0}^{c-1} \tilde{Q}_{jo} \cdot p(y = j \mid x),$$

where  $\tilde{Q}_{jo} = p(\tilde{y} = o \mid y = j)$  and  $o = tm + k$ . Then we'll have target-location-based risk:

$$\mathcal{R}_{\text{TL}}(h) = \mathbb{E}_{p(x, \tilde{y})} \ell(\tilde{q}(x), \tilde{y}).$$



# Theorem

**Assumption.**  $\forall Y \in \mathcal{S}(\tau), p(Y) = 1/|\mathcal{S}(\tau)|$ .

**Theorem 1.**  $\mathcal{R}_L(h) \leq \mathcal{R}_{\text{NeSy}}(h) + C(\mathcal{S}(\tau))$ .

**Theorem 2.**  $\mathcal{R}_{\text{TL}}(h) \leq \mathcal{R}_{\text{NeSy}}(h) + C(p(\tau))$ .

**Theorem 3.** If  $\tilde{Q}$  has full row rank, then  $h_{\text{TL}}^* = \operatorname{argmin}_h \mathcal{R}_{\text{TL}}(h)$  recovers  $h^* = \operatorname{argmin}_h \mathcal{R}(h)$ , i.e.,  $h_{\text{TL}}^* = h^*$ .

**Corollary 1.**  $h_L^* = h^*$ .

# Example

## Facts

```
conj([0,0,0]).
conj([0,1,0]).
conj([1,0,0]).
conj([1,1,1]).
```

## Knowledge Base

```
conj([Y0,Y1,Y2]) ← Y2 is Y0 ∧ Y1.
```

## Raw Data

```
0 0 0  0 1 0  1 0 0  1 1 1
```

## Facts

```
conj0([0,0]).
conj0([0,1]).
conj0([1,0]).
conj1([1,1]).
```

## Knowledge Base

```
conj([Y0,Y1,Y2]) ← See Figure 2.
conj0([Y0,Y1]) ← conj([Y0,Y1,0]).
conj1([Y0,Y1]) ← conj([Y0,Y1,1]).
```

$$Q = \begin{pmatrix} 2/7 & 2/7 & 3/7 \\ 2/5 & 2/5 & 1/5 \end{pmatrix}$$

$$Q = \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix}$$

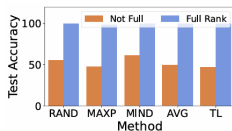
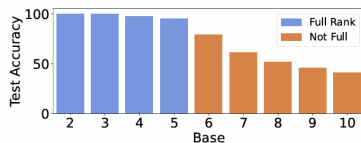
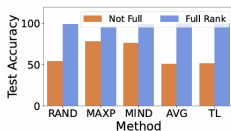
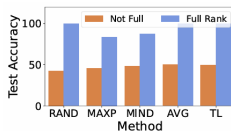
$$\tilde{Q} = \begin{pmatrix} 1/2 & 1/2 & 0 & 0 \\ 1/4 & 1/4 & 1/4 & 1/4 \end{pmatrix}$$

# Experiment

TASK	METHOD	MNIST	EMNIST	USPS	KUZUSHIJI	FASHION
ConjEq	RAND	99.91 $\pm$ 0.06	99.65 $\pm$ 0.04	99.33 $\pm$ 0.16	97.82 $\pm$ 0.35	98.40 $\pm$ 0.11
	MAXP	99.94 $\pm$ 0.04	99.82 $\pm$ 0.03	99.20 $\pm$ 0.00	98.80 $\pm$ 0.16	99.39 $\pm$ 0.12
	MIND	99.91 $\pm$ 0.08	99.84 $\pm$ 0.07	99.14 $\pm$ 0.17	98.91 $\pm$ 0.17	98.84 $\pm$ 0.19
	AVG	99.85 $\pm$ 0.10	99.80 $\pm$ 0.07	99.30 $\pm$ 0.17	98.34 $\pm$ 0.16	98.62 $\pm$ 0.21
	TL	99.92 $\pm$ 0.05	99.82 $\pm$ 0.06	99.25 $\pm$ 0.08	98.53 $\pm$ 0.26	98.77 $\pm$ 0.06
Conjunction	RAND	99.91 $\pm$ 0.06	99.86 $\pm$ 0.04	99.30 $\pm$ 0.13	98.79 $\pm$ 0.13	99.00 $\pm$ 0.27
	MAXP	99.93 $\pm$ 0.04	99.81 $\pm$ 0.02	99.20 $\pm$ 0.00	98.62 $\pm$ 0.15	99.05 $\pm$ 0.09
	MIND	99.94 $\pm$ 0.02	99.79 $\pm$ 0.02	99.20 $\pm$ 0.00	98.74 $\pm$ 0.10	99.08 $\pm$ 0.10
	AVG	99.94 $\pm$ 0.02	99.85 $\pm$ 0.03	99.30 $\pm$ 0.13	98.68 $\pm$ 0.33	99.23 $\pm$ 0.13
	TL	99.94 $\pm$ 0.02	99.83 $\pm$ 0.04	99.20 $\pm$ 0.00	98.87 $\pm$ 0.18	99.30 $\pm$ 0.04
Addition	RAND	92.01 $\pm$ 0.93	92.94 $\pm$ 1.45	90.96 $\pm$ 1.04	73.18 $\pm$ 0.71	79.08 $\pm$ 2.61
	MAXP	96.40 $\pm$ 4.04	95.09 $\pm$ 5.20	94.29 $\pm$ 0.27	90.00 $\pm$ 0.27	87.34 $\pm$ 2.93
	MIND	98.32 $\pm$ 0.04	98.61 $\pm$ 0.06	94.61 $\pm$ 0.17	90.85 $\pm$ 0.26	88.40 $\pm$ 0.62
	AVG	94.90 $\pm$ 0.39	95.71 $\pm$ 0.42	93.22 $\pm$ 0.30	80.94 $\pm$ 0.62	84.43 $\pm$ 0.92
	TL	98.00 $\pm$ 0.14	98.41 $\pm$ 0.05	94.68 $\pm$ 0.20	90.04 $\pm$ 0.32	88.38 $\pm$ 0.25
HED	RAND	99.89 $\pm$ 0.02	99.71 $\pm$ 0.12	99.25 $\pm$ 0.23	97.68 $\pm$ 0.70	98.43 $\pm$ 0.55
	MAXP	99.90 $\pm$ 0.02	99.77 $\pm$ 0.02	99.23 $\pm$ 0.05	98.55 $\pm$ 0.08	99.33 $\pm$ 0.10
	MIND	99.87 $\pm$ 0.07	99.77 $\pm$ 0.02	99.21 $\pm$ 0.00	98.61 $\pm$ 0.21	99.32 $\pm$ 0.10
	AVG	99.60 $\pm$ 0.09	99.38 $\pm$ 0.21	99.32 $\pm$ 0.14	96.16 $\pm$ 1.22	98.46 $\pm$ 0.33
	TL	99.90 $\pm$ 0.02	99.77 $\pm$ 0.04	99.21 $\pm$ 0.00	98.50 $\pm$ 0.16	99.21 $\pm$ 0.06

Table 1: Test accuracy (%) of each method using MLP on benchmark datasets and tasks.

# Experiment

(a) DNF,  $m = 3$ (b) DNF,  $m = 4$ (c) DNF,  $m = 5$

# My Works

What if there are less columns than rows in  $Q$ ?

$$\text{SimilarityScore}(h) = \text{Inter}(h) - \text{Intra}(h)$$

$$\text{Inter}(h) = \sum_{i_1 \neq i_2, j_1 \neq j_2} \omega_{i_1 j_1} \omega_{i_2 j_2} \text{Dis}(x_{i_1}, x_{i_2})$$

$$\text{Inter}(h) = \sum_{i_1 \neq i_2, j} \omega_{i_1 j} \omega_{i_2 j} \text{Dis}(x_{i_1}, x_{i_2})$$

If any two rows of  $Q$  are linearly independent. By maximising the SimilarityScore after minimising  $\mathcal{R}_L$ , we may still obtain a good classifier.

**Thm 0.** 令 SimilarityScore 最小等价于令

$S_s = \sum_{i_1, i_2} \text{Dis}(x_{i_1}, x_{i_2}) \cdot \omega_{i_1} \omega_{i_2}$  最小.

**Thm 1.** 若距离的定义足够好, 即任意相同类图像间距离均小于不同类图像间距离, 并且训练集中, 属于每一类的图像一样多时, 那么当  $h$  将不同类图像完全区分开时,  $S_s$  取最小值。

**Thm 2.** 令  $\omega = \tilde{\omega}$  时,  $S_s$  可被  $S_s$  的最小值控制。

Thank you

Thank you!