

# **Algebraic Logic for Deep Learning**

by

YKY

A Thesis Submitted to  
The Hong Kong University of Science and Technology  
in Partial Fulfilment of the Requirements for  
the Degree of Master of Philosophy  
in Applied Mathematics

August 2024, Hong Kong



## **Authorization**

I hereby declare that I am the sole author of the thesis.

I authorize the Hong Kong University of Science and Technology to lend this thesis to other institutions or individuals for the purpose of scholarly research.

I further authorize the Hong Kong University of Science and Technology to reproduce the thesis by photocopying or by other means, in total or in part, at the request of other institutions or individuals for the purpose of scholarly research.

---

YKY

13 August 2024



# Algebraic Logic for Deep Learning

by

YKY

This is to certify that I have examined the above MPhil thesis  
and have found that it is complete and satisfactory in all respects,  
and that any and all revisions required by  
the thesis examination committee have been made.

---

Prof. Who, Thesis Supervisor

---

Prof. Who Else, Thesis Co-supervisor

---

Prof. Ikari Someone, Head of Department

Department of Maths

13 August 2024



## **ACKNOWLEDGEMENTS**

Thank you, all the Evangelion.





# TABLE OF CONTENTS

Title Page	i
Authorization	iii
Signature Page	v
Acknowledgements	vii
Table of Contents	ix
List of Figures	xiii
List of Tables	xv
List of Algorithms	xvii
Abstract	xix
Chapter 0    Introduction	1
Chapter 1    Background: The mind as a dynamical system	3
1.1    The set-up	3
Chapter 2    Background: Categorical logic	5
2.1    Topos and internal language	5
Chapter 3    Background: Algebraic logic	7
3.1    Paul Halmos' algebraic logic	7
3.2    Yuri Manin and Russians	7
3.3    Term rewriting and all that	7
Chapter 4    Design of algorithm	9
4.1    From abstract algebraic logic to concrete computations	9
4.2    What does it mean to train the AI?	9

4.3	“Geometric” logic inference algorithm	12
4.3.1	How to determine if a rule is satisfied	12
4.3.2	How to handle variables in rules	13
4.4	Computer representation of rules	14
4.5	Rules recommender	15
4.5.1	Differentiability	15
4.6	Interestingness	17
4.7	Combining logic and reinforcement learning	17
4.7.1	Logical policy function	17
4.8	Basics of DQN (Deep Q Learning)	18
Chapter 5	Combining RL and Auto-regression	21
Chapter 6	Tic Tac Toe experiment	23
6.1	Representation of states and rules	23
Chapter 7	Conclusions	25
Appendix A	List of Publications	27
Appendix B	FYTGS Requirements	29
B.1	Components	29
B.1.1	Order	29
B.1.2	Authorization page	29
B.1.3	Signature page	30
B.1.4	Acknowledgments	30
B.1.5	Abstract	30
B.1.6	Bibliography	30
B.2	Language, Style and Format	30
B.2.1	Language	30
B.2.2	Pagination	31
B.2.3	Format	31
B.2.4	Footnotes	31
B.2.5	Appendices	32
B.2.6	Figures, Tables and Illustrations	32
B.2.7	Photographs/Images	32





## **LIST OF FIGURES**



## **LIST OF TABLES**





## **LIST OF ALGORITHMS**



# **Algebraic Logic for Deep Learning**

by YKY

Department of Maths

The Hong Kong University of Science and Technology

Abstract

Some text.



# **CHAPTER 0**

## **INTRODUCTION**

Good luck.



# CHAPTER 1

## BACKGROUND: THE MIND AS A DYNAMICAL SYSTEM

### 1.1 The set-up

The set of equations  $F$  defines an algebraic set = **the world**:

$$F(x) = 0. \quad (1.1)$$

The objective of an intelligent agent is to learn  $F$ .

We have the function  $f$  performing **prediction** of the immediate future:

$$\boxed{\text{current state}} \quad x_t \xrightarrow{F} x_{t+1} \quad \boxed{\text{next state}}. \quad (1.2)$$

In an infinitesimal sense, we can see  $f$  as a **differential equation** describing the **world trajectory**:

$$\dot{x} = f(x). \quad (1.3)$$

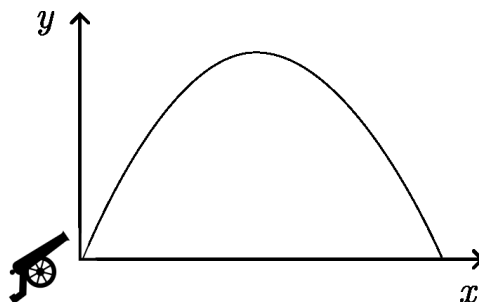
So  $F$  is the **solution** to this differential equation.

It seems that  $F$  and  $f$  are more or less equivalent ways to describe the world.

Logic can be turned into some form of algebra, and this algebra can be used to express either  $F$  or  $f$ . Perhaps both ways are feasible, or even mixing the two.

What does it mean to use logic to express  $F$  or  $f$ ?

Go back to a physics example, the parabolic trajectory of a canon ball:



(1.4)

The parabola is given by the quadratic equation from high school:

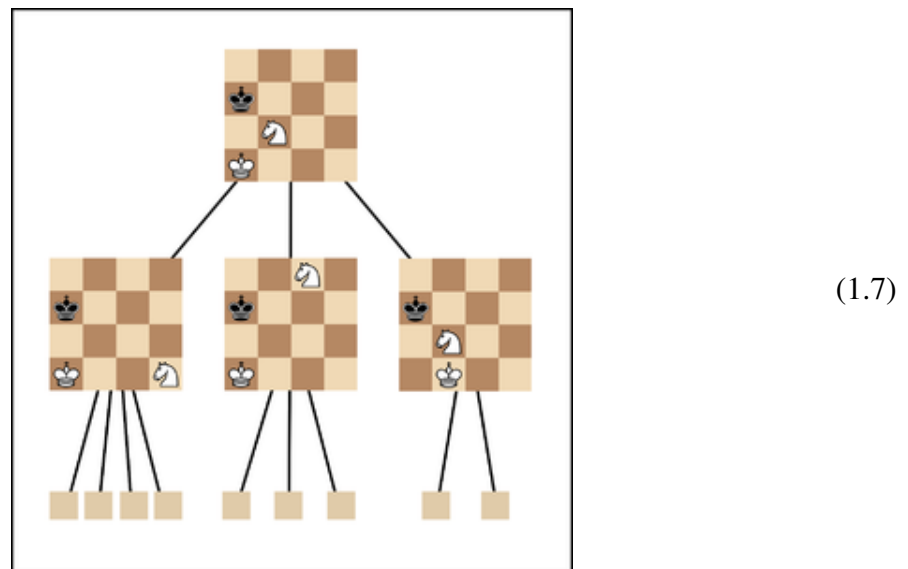
$$F(x) = 0 \quad F(x) = ax^2 + bx + c \quad (1.5)$$

but the trajectory can also be described by the physics equation parameterized by time  $t$ :

$$\dot{\mathbf{x}} = f(\mathbf{x}) = (v_x, v_y) = (v_x^0, -gt + v_y^0). \quad (1.6)$$

This parametric form is not unique. For example, another way is for the point  $\mathbf{x}$  to move with uniform speed along the trajectory.

Note that the trajectory above is qualitatively the same as a “thought trajectory” in cognitive space. One intuitive way to visualize cognitive trajectories is via the example of a chess game tree (where the opponent’s move is regarded as how the “world” reacts to the agent’s actions):



$f(x)$  is the functional form we want, as it contains information about the “interestingness” of deduced conclusions.

Are our equations in Table 4.5 describing  $F$  or  $f$ ?

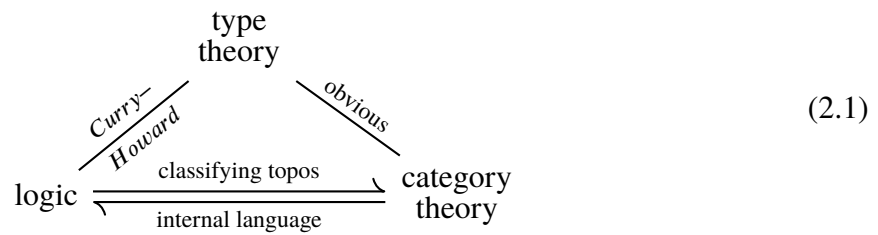
An intuitive idea: state = facts = grounded equations, rules = quantified equations. So the equations are modeling  $f$ . This exposed a problem of classical AI that I have not paid too much attention to: the selection of interesting conclusions. It’s hard to **enumerate conclusions**, let alone to rate their interestingness.



# CHAPTER 2

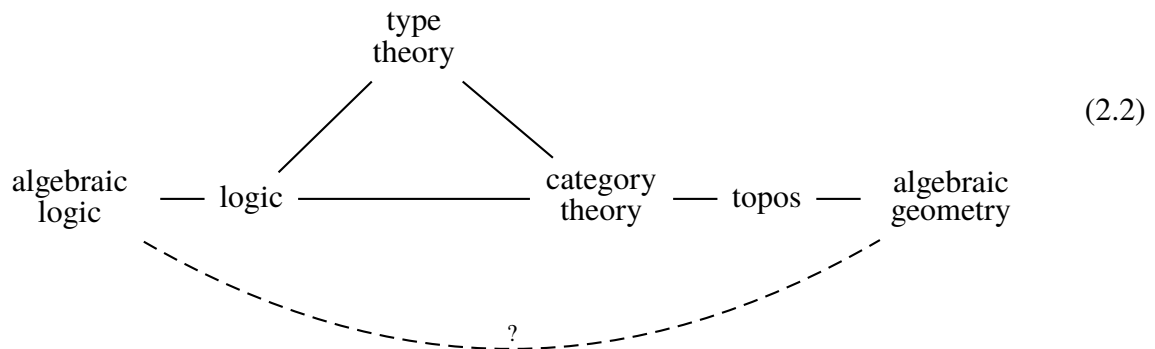
## BACKGROUND: CATEGORICAL LOGIC

Lambek posited this “trinity”:



where the double arrows at the base can be understood thusly:

I extended some nodes to better see their relations:



I am curious if the two algebras on the left and right are identical?

### 2.1 Topos and internal language



## CHAPTER 3

### BACKGROUND: ALGEBRAIC LOGIC

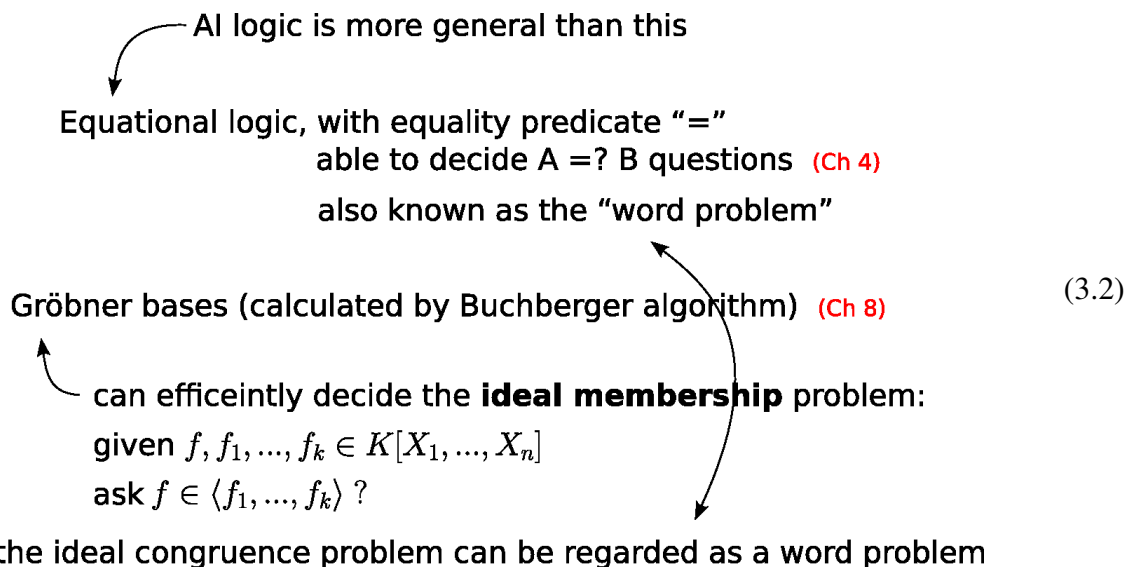
#### 3.1 Paul Halmos' algebraic logic

Every Boolean algebra  $\mathbb{A}$  is isomorphic to the set of all continuous functions from  $X$  into  $\mathbb{O}$ , where  $X$  is the dual space of the algebra  $\mathbb{A}$ , and  $\mathbb{O}$  is the Boolean algebra with 2 elements. If there is a homomorphism  $f$  between Boolean algebras  $\mathbb{A} \rightarrow \mathbb{B}$  then there is a dual morphism  $f^*$  between their dual spaces  $Y \rightarrow X$ :

$$\begin{array}{ccc}
 \mathbb{A} & \xrightarrow{f} & \mathbb{B} \\
 \cong & & \cong \\
 \overline{X} & \xleftarrow{f^*} & \overline{Y} \\
 \downarrow & & \downarrow \\
 \mathbb{O} & & \mathbb{O}
 \end{array} \tag{3.1}$$

#### 3.2 Yuri Manin and Russians

#### 3.3 Term rewriting and all that



- What is the **word problem**?

is defined for an equational theory  $E$ .

is the problem of deciding whether  $s = t$

- why is Gröbner basis equivalent to the word problem

to ask ideal congruence  $f = ?g$  means  $f - g \in ?J$

which is ideal membership problem

a polynomial can be regarded as a rewrite rule

because  $f = 0$ , we can take the “largest monomial” in  $f$  as the LHS, and the rest of  $f$  as RHS.

In other words: ideal = set of rules

We ask if a polynomial can be rewritten by the set of rules to another form.

This is similar to logic deduction.

- Here an important question is: polynomial reduction seems unable to handle **logic variables**, it seems only capable of **simple symbolic rewriting**.
- logic is equivalent to what form of polynomials?  
taking the cue that Huet’s higher-order unification = Buchberger algorithm, ...

# CHAPTER 4

## DESIGN OF ALGORITHM

The following table depicts the main correspondences relevant to our research:

<b>LOGIC</b>	<b>facts</b> human(socrates)	<b>rules</b> $\forall x.\text{human}(x) \rightarrow \text{mortal}(x)$
<b>ALGEBRA</b>	<b>element</b> $p \in \mathbb{A}$	<b>element</b> $(p \rightarrow q) \in \mathbb{A}$
<b>WORLD</b>	<b>states</b> $x_t$	<b>state transitions</b> $\overset{F}{x_t \mapsto x_{t+1}}$

(4.1)

The relation between LOGIC and WORLD has been elucidated quite thoroughly in the AI literature. Note that the state  $x_t$  is made up of a set of facts (logic propositions). A single step of logic inference results in a new conclusion  $\delta x$  which is *added* (as a set element) to the current state  $x_t$  to form a new state  $x_{t+1}$ . Here  $t$  refers to “mental time” which does not necessarily coincide with real time.

### 4.1 From abstract algebraic logic to concrete computations

There are two main routes to make abstract algebraic logic concrete:

- Find **matrix representations** of the logical algebra
- Implement the logical algebra as the commutative algebra of (classical) **polynomials**

### 4.2 What does it mean to train the AI?

From the previous section,

$$F(x) = 0 \quad \text{is the solution to} \quad \dot{x} = f(x) \tag{4.2}$$

and the two descriptions (by  $F$  or by  $f$ ) are equivalent.

The discrete version of  $f$  is  $\mathcal{f}$ :

$$\mathcal{f}(x_t) = x_{t+1} - x_t = \delta x \quad (4.3)$$

The sensory data from the AI are a set of “world” points  $\{x_i\}$  and we require either:

$$F(x_t) = 0 \quad \text{or} \quad \mathcal{f}(x_t) = \delta x = x_{t+1} - x_t \quad (4.4)$$

and  $F$  or  $f$  can be trained by gradient descent to eliminate errors in the above conditions (equations).

- the  $x_t$ 's are represented as **logic facts**
- $F$  or  $f$  is represented as **logic rules**

and we need to **evaluate**  $F(x_t)$  or  $f(x_t)$ .

Let's do some examples:

Logic formula	Algebraic form
human(socrates)	$h(s) = 1$
human(socrates) $\wedge$ human(plato)	$h(s) \cdot h(p) = 1$
human(socrates) $\rightarrow$ mortal(socrates)	$1 + h(s) + h(s) \cdot m(s) = 1$
$\forall x. \text{human}(x)$	$h(x)$ is a propositional function $\forall_x h(x)$ is a constant function mapping to 1 or 0
$\forall x. \text{human}(x) \rightarrow \text{mortal}(x)$	$\forall_x (1 + h(x) + h(x) \cdot m(x))$ is a constant function mapping to 1 or 0
$\forall x, y, z. \text{father}(x, y) \wedge \text{father}(y, z) \rightarrow$ grandfather(x, z)	$\forall_x \forall_y \forall_z (1 + f(x, y) \cdot f(y, z) + f(x, y) \cdot f(y, z) \cdot g(x, z))$ $\mapsto 0 \text{ or } 1$
<b>general Horn formula:</b> $\forall_{x...} P \wedge Q \wedge R... \rightarrow Z$	$\forall_{x...} (1 + P \cdot Q \cdot R... + P \cdot Q \cdot R... \cdot Z)$ $\mapsto 0 \text{ or } 1$

(4.5)

Now imagine there are millions of such rules. Number of predicates obviously increases.

Does each equation require new variables, or can variables be re-used? Seems yes, can be re-used.

The **loss function** would be the sum of squared errors over all equations:

$$\mathcal{L} = \sum_{\text{eqns}} \epsilon^2 = \sum_i (\phi_i(x...) - 1)^2. \quad (4.6)$$

Learning means to perform the **gradient descent** via  $\nabla_{\Phi} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \Phi}$  where  $\Phi$  is the set of parameters for the equations.

Potential problems:

- How to represent the set of equations efficiently? Matrix of coefficients seems wasteful.
- We have lost the “**deepness**” of deep learning, but there is recent research showing that **shallow learning** may work well too.
- Need to iterate logical inference multiple times using the same set of equations.

How are new conclusions added to the state? What is the state? State = set of facts = set of **grounded** equations.

Inference: how to get from current state to next state? Big problem!! New ground facts have to be read off from satisfaction of all equations. Rather intractable...

See Chapter 1.

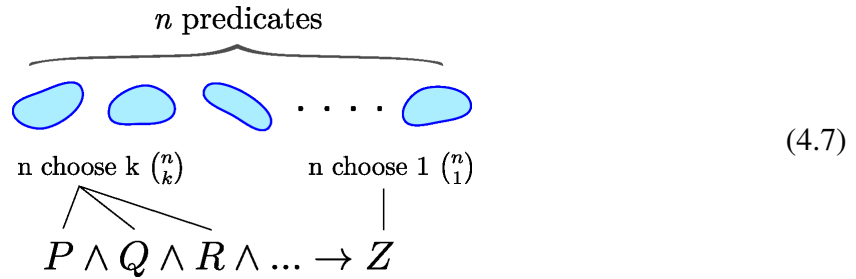
How to enumerate conclusions? The loss function (4.6) can be trained on any data with correlations. But what we have is data in the form of time series. Need extra measures to ensure that equations only model  $x_t \rightarrow x_{t+\Delta t}$ . Now how to enumerate conclusions? From current state, iterate over all equations to generate new states. This is getting very close to the classical logic-based AI inference algorithm.

The **interestingness function** gives a probability distribution over conclusions given the current “context” (which we identify as the current state  $x_t$ ):  $\text{Intg}(\delta x) = \mathbb{P}(\delta x | x_t)$ . This function has to be learned. It has an equivariant structure due to the state as a set of propositions with permutation invariance.

One more efficiency problem: iterating through all equations is inefficient, which brings back the necessity of the classical **rete algorithm**: instead of matching rules against the state, we should match the current state against rules. In other words, instead of  $\delta x := \bigcup_i \text{rule}_i(x)$ , perform  $\delta x := \text{compiled-rules}(\Delta x)$ , where  $\Delta x$  is the change of the current state from the previous state, and  $\delta x$  is the change from the current state to the next state.

But we may also avoid the rete algorithm by stacking logical equations into **layers**, thus getting an efficiency advantage similar to deep learning. A “single” step of logic inference would mean going through multiple layers of logic rules (equations).

### 4.3 “Geometric” logic inference algorithm



- What is a logic fact within a state?
- How does a rule generate a single new fact?

A fact consists of a point (in subject space) and a predicate that contains it. The point itself does not suffice because it can belong to various predicates.

#### 4.3.1 How to determine if a rule is satisfied

To apply a rule, each **atomic term** in the rule has to be satisfied. For each predicate  $Q()$ , this is verified by testing if we have any points among the facts contained in  $Q$ .

If we have `father(john, pete)` as a fact then we certainly can satisfy `father(x, y)`. But we already have the point (john, pete) which may satisfy other predicates  $Q(x, y)$ . So our method is slightly more permissive (and thus more powerful) in rules matching.

How does the rule’s RHS generate a new fact? It should also be a (point, predicate) pair. The point has to **match** the premise. How could this be ensured?

Secondly, the output predicate may not cover the point.

The matching process: Syntactically, we look at each literal in the rule and see if any fact unifies with the literal. Geometrically, it means taking a point and checking if it lies inside a predicate. If the fact is a (point, predicate) pair then it is given that the point belongs to the predicate, so it is not necessary to check for membership. The result is simply taking the point when the predicates match. But we still need to keep track of which variables the point coordinates are binding to.

The matching of the second literal will also return a point, but its coordinates would be bound to different dimensions.

The binding of coordinates would be the basis of verifying the rule LHS.

RHS: The bound variables (in various dimensions) need to be projected to the output space. It may or may not be covered by the output predicate. If not, the rule does *not* apply. This is in accord with the principle that rules should not change during inference.



### 4.3.2 How to handle variables in rules

( Besides  $(p, \gamma)$  we seem to need another parameter, the variable number specifying which of  $\{v_1, \dots, v_n\}$  is referenced in a rule's argument. This seems to be a parameter independent of  $(p, \gamma)$ . But it seems irrelevant if there is no cylindrification? )

Again drawing inspiration from Halmos' algebraic logic. A monadic logic deals with predicates as  $X \rightarrow \Omega$ , whereas in polyadic logic one has  $X^I \rightarrow \Omega$  where  $I$  is an index set. The elements of  $I$  are variables even though the elements themselves do not really change. This creates  $I$  copies of the Subject Space  $X$ <sup>1</sup>. For example if  $I = \{1, 2, 3, 4\}$  would provide 4 variables or “slots” for a rule to use. A rule such as  $\text{father}(X, Y) \wedge \text{father}(Y, Z) \rightarrow \text{grandfather}(X, Z)$  requires 3 variables:  $X, Y, Z$ .

Analyzing this example in more details:

$$\text{father}(X_1, X_2) \wedge \text{father}(X_2, X_3) \rightarrow \text{grand-father}(X_1, X_3) \quad (4.8)$$

where I have used the symbol  $X$  with subscripts to denote variables, as a reminder that each variable is really a **copy** of our Subject Space  $X$ .

It is perhaps illuminating to rewrite the logic rule this way:

$$\text{father}(X^I) \wedge \text{father}(X^I) \rightarrow \text{grand-father}(X^I) \quad (4.9)$$

where  $I$  is at least  $\{1, 2, 3\}$ .

In my mind, the mental picture is like this:

$$\text{father}(\text{||||}) \wedge \text{father}(\text{|||}) \rightarrow \text{grand-father}(\text{|||}) \quad (4.10)$$

where the blue bars represent fuzzy-values, some kind of relative “proportion” of variables, such that the overall construction would be differentiable. ( Note that  $X_1$  does not necessarily map to a leftmost dimension, it could map to any slot further to the right. )

The solution I found is quite natural: each **argument** of a predicate would be associated with a **softmax** over the index set  $I$ , which selects 1 element out of  $I$ . The softmax assigns, to each  $i \in I$ , a weight  $w_i \in \mathbb{R}$ . Then we can make copies of  $X$  into  $X^I$  as  $\{w_1 X_1, \dots, w_n X_n\}$ . But this method assumes that the Subject Space  $X$  has **vector space structure** and is amenable to scalar multiplication. As I have argued in a blog article, the space of “word embedding” is actually a

<sup>1</sup>Note that  $(I \rightarrow X) \rightarrow \Omega$  is not isomorphic to  $I \rightarrow (X \rightarrow \Omega)$ .

**metric space** rather than vector space.

A better solution is to associate the weight to that argument of the predicate such that the resulting **true value** of that predicate would be **weakened**. The truth value of a predicate  $P$  is evaluated by sending a point  $(a, b)$  as input to the neural network representing  $P$ . Its output is a real number normalized  $\in [0, 1]$  and interpreted as a fuzzy truth-value.

## 4.4 Computer representation of rules

Each logic rule may vary in the following ways:

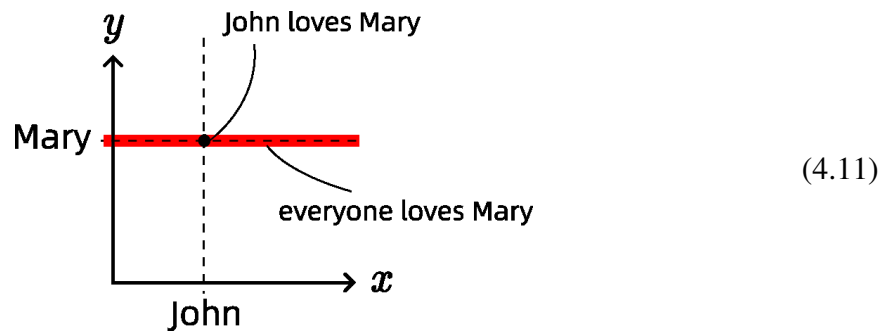
- number of literals and their polarity
- number of arguments
- each argument may be a logic variable or logic constant
- which specific logic variable or constant.

#1 and #2 can be fixed.

For #3 and #4, each variable or constant may be represented by a pair  $(p, \gamma)$  where  $p$  is a point in the Subject Space, and  $\gamma$  is the **cylindrification factor** representing the degree to which this argument is like a point or like a “for all” variable. But this is problematic as the dimension of the Subject Space  $X$  is not the same as the dimension of the variable space  $\{X_1, X_2, \dots, X_n\}$ .

In fact, each predicate has as its arguments  $n$  copies of the Subject Space, ie,  $X^n$ . One or more of the copies may be **cylindrified**.

There are two senses of the word “cylindrify.” Think for example the change from (John, Mary) to  $(x, \text{Mary})$  as in “everyone loves Mary.” In the first sense, the point John on the  $x$ -coordinate becomes “everyone,” ie, the entire  $x$ -dimension becomes a cylinder. In the second sense, a cylinder appears on the  $y$ -axis at the position of Mary:



In the above example, we have two copies of the Subject Space  $X$ , each is 1-dimensional.

The dimension  $n$  is not just the number of variables appearing in one predicate in a rule, but the total number of variables appearing in a rule.

---

We suggest the following values (for a general human-level intelligent agent):

Variable	Symbol	Suggested values
# of rules	$M$	millions up
# of literals per rule	$K$	2-10, typical 3,4
# of arguments per rule	$n$	2, 3, 4
Subject Space dimension	$\dim X$	100-1000, typical 512, 768

(4.12)

Is everything differentiable?

## 4.5 Rules recommender

The rules recommender would be a set function:

$$R_u : \{\text{current state}\} \rightarrow \{\text{set of rules}\} \quad (4.13)$$

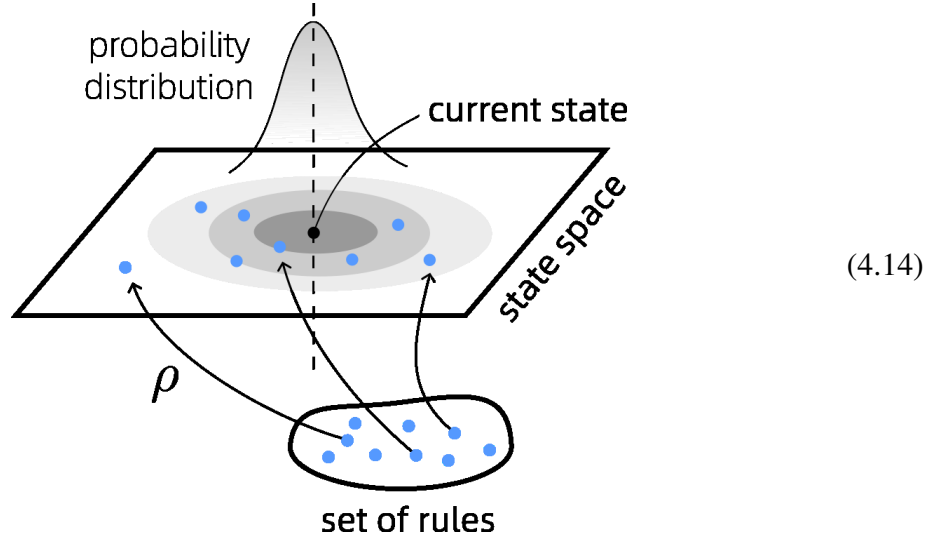
which has equivariant structure on both its input and output. This suggests the **Transformer** architecture is suitable for learning this function.

But this is clearly **non-differentiable**!

### 4.5.1 Differentiability

In binary logic a rule either applies or does not apply. To make the entire set of rules differentiable, rule application must be “graded”.

So we propose a **probabilistic** rule-matching mechanism:



Rules are mapped by a function  $\rho$  to the state space, such that rules are located close to the states they are most likely to match (this is a multi-objective optimization problem). Given a current state  $x$ , denote the ball around it as  $B(x, d)$ . The rules we predict most likely to match  $x$  are given by  $\rho^{-1}(B(x, d))$  for some distance  $d$ . We can assign a Gaussian probability distribution centered around  $x$  over such rules, so that rules are selected **probabilistically** for update.

It seems that doing so would not affect the convergence (of the learning algorithm) of the rules, but it would be much more efficient as the probability distribution is *more* concentrated around  $x$ , so that large numbers of rules need *not* be evaluated (for each state, at each time step).

$$x_{t+1} = F(x_t; \Theta) \quad (4.15)$$

where  $\Theta$  are all the parameters that determine the rules,  $F$  is the total function aggregating all the rules. If we vary any one rule parameter, does it change abruptly? The parameters  $\Theta$  include the rules-sorting function  $\rho$ .

For gradient descent we need to calculate  $\nabla_{\Theta} \mathcal{L}$  where  $\Theta$  are all the parameters of rules, and the loss function is defined as the sum of errors over all rules,  $\mathcal{L} = \sum_{\text{all rules}} e^2$ . Each error usually is given by the difference as compared against a reference value (supervised learning), but in our case such a value is unavailable, instead the objective comes from reinforcement learning, ie, the Hamilton-Jacobi-Bellman equation:  $J(x_{t+1}) = \max_a [\gamma J(x_t) + R(x_t, a)]$ , where  $a$  means **action**. In the logic-based setting, an action means applying a logic rule to change the state.

Policy function:

$$\pi : X \times A \rightarrow [0, 1] \quad (4.16)$$

By the **Policy Gradient Theorem**, calculation of  $\nabla_{\Theta} J$  translates to calculating  $\nabla_{\Theta} \pi$ .

In reinforcement learning, in particular the **Policy Gradient** algorithm, we need to calculate the gradient  $\nabla_{\phi} \mathcal{L}(\phi)$  of a loss function of the form:

$$\mathcal{L}(\phi) = \mathbb{E}_{z \sim q_{\phi}(z)}[f(z)]. \quad (4.17)$$

The expectation gives an integral  $\nabla_{\phi} \mathcal{L}(\phi) = \nabla_{\phi} \int dz q_{\phi}(z) f(z)$  which removes the “randomness” and evaluates to a real number. This allows to be handled by traditional differential calculus.

At each state, the application of all rules results in a list of all available actions. The differentiability problem may arise from probabilistic rule-matching.

## 4.6 Interestingness

This can be implemented implicitly by making the inference algorithm output a probability distribution over all deduced conclusions and then picking the most probable one.

## 4.7 Combining logic and reinforcement learning

### 4.7.1 Logical policy function

The policy function  $\pi : X \times A \rightarrow [0, 1]$  should be implemented with logic structure. The current state would be matched against rules selected by the rules recommender, and then each rule would be applied, resulting in conclusions  $Z_i$  each associated with a probabilistic strength  $\mathbb{P}(Z_i)$ . These are the available actions and their probabilities as given by the policy.

---

This section concerns how to establish the “link” from logic to reinforcement learning.

For reinforcement learning, I find it easier to consider Q-learning, ie, learning the utility function  $Q(s, a)$ , where  $s$  = current state and  $a$  = action taken in current state.

During forward inference, each logic rule of the form (also known as Horn form):

$$P \wedge Q \wedge R \wedge \dots \rightarrow Z \quad (4.18)$$

yields a conclusion  $Z$ . But the “truths” of the premises  $P, Q, R, ..$  are not binary but fuzzy, because the rules have to be differentiable (this can be implemented by softmax). Thus each conclusion  $Z$  is also fuzzy. The application of all the rules yields a set of conclusions  $\{Z_i\}$ , with their associated fuzzy truth values, which we can normalize as a probability distribution over all available next actions. This can be regarded as the **policy function**  $\pi(a|s)$ , which gives the probability of an action given the current state,  $\mathbb{P}(a|s)$ .

In Q-learning we need to learn the values  $Q(s, a)$ . Given the reward  $R$ ,  $Q(s, a)$  satisfies the Bellman equation:

$$Q(s, a) += \eta (R + \gamma \max_{a'} Q(s', a') - Q(s, a)) \quad (4.19)$$

where  $\eta$  is the learning rate,  $\gamma$  is the discount rate of values.

The policy function  $\pi(a|s)$  is not exactly the same as  $Q(s, a)$ , but there is a trick that can relate them, that is the basis of **Soft-Q learning** (Soft Actor-Critic, SAC):

$$\pi(a|s) \propto \exp Q(s, a) \quad (4.20)$$

The situation can be illustrated by the following figure (not a commutative diagram):

$$\begin{array}{ccc} & & \text{Bellman} \\ & & \text{update} \\ & & \downarrow \\ \text{logic} & \longrightarrow & \pi(a|s) \xrightleftharpoons[\text{exp}]{\text{log}} Q(s, a) \end{array} \quad (4.21)$$

Logic rules give us  $\pi(a|s)$ , which in turn gives us  $Q(s, a)$ , but  $Q(s, a)$  is also determined by Bellman update.

My solution is: Logic outputs a finite set of actions, to each action is associated a weight, which we identify as  $Q$ . Use Bellman update to train these  $Q$  values. To choose an action, calculate  $\pi(a|s)$  as suggested by the Soft-Q method. This has the effect of **maximum-entropy** exploration.

## 4.8 Basics of DQN (Deep Q Learning)

The **Bellman optimality condition** says that:

$$Q_t^* = \max_a \{ R + \gamma Q_{t+1}^* \} \quad (4.22)$$

where  $()^*$  denotes optimal values.

In the **Bellman update**,

$$Q_t += \eta \cdot \text{TD-error} \quad (4.23)$$

where TD = **temporal difference** is defined as:

$$\text{TD-error} = \overbrace{R + \gamma \max_{a'} Q(s_{t+1}, a')}^{\text{ideal value}} - \overbrace{Q(s, a)}^{\text{current value}} . \quad (4.24)$$

So obviously the Bellman update causes the current value of  $Q$  to **converge** to that of the ideal value  $Q^*$ .

Remember that, in the simplest **Q-table** method, we apply the Bellman update over Q-table entries. The **DQN** approach differs in that, instead of a Q-table, we use a deep neural network in its place. So the updating of the Q-table now becomes updating the DQN via **gradient descent**.





# **CHAPTER 5**

## **COMBINING RL AND AUTO-REGRESSION**



# CHAPTER 6

## TIC TAC TOE EXPERIMENT

### 6.1 Representation of states and rules

At any time the state is a set of facts, ie, pairs of (point  $\in$  predicate).

There will be  $K$  predicates and  $M$  rules.

Each rule is a conjunction of all predicates. If  $K$  is large, the rules would be cumbersome.

Having a large number of rules,  $M$ , seems not to have a deleterious effect, if the rules recommender is good at its job.

Each literal in the rule may be negated, how to handle this?

Each literal contains a predicate and its arguments, which can be constants or variables.

It seems that genetic algorithms would be best suited for this kind of search for rules... or unless the rules are represented in such a way that they can continuously vary in a differentiable manifold.

For TicTacToe, let's say number of rules = 50.



# **CHAPTER 7**

## **CONCLUSIONS**

Some conclusion text.



**APPENDIX A**

**LIST OF PUBLICATIONS**





# APPENDIX B

## FYTGS REQUIREMENTS

The requirements are from the RPG Handbook.

### B.1 Components

#### B.1.1 Order

A thesis should contain the following parts in the order shown:

1. Title page, containing in this order:
  - a. Thesis title
  - b. Full name of the candidate
  - c. Degree for which the thesis is submitted
  - d. Name of the University, *i.e.* The Hong Kong University of Science and Technology
  - e. Month and year of submission
2. Authorization page
3. Signature page
4. Acknowledgments
5. Table of contents
6. Lists of figures and tables
7. Abstract ( $\leq 300$  words.)
8. Thesis body
9. Bibliography
10. Appendices and other addenda, if any.

#### B.1.2 Authorization page

On this page, students authorize the University to lend or reproduce the thesis.

1. The copyright of the thesis as a literary work vests in its author (the student).

2. The authorization gives HKUST Library a non-exclusive right to make it available for scholarly research.

### **B.1.3 Signature page**

This page provides signatures of the thesis supervisor(s) and Department Head confirming that the thesis is satisfactory.

### **B.1.4 Acknowledgments**

The student is required to declare, in this section, the extent to which assistance has been given by his/her faculty and staff, fellow students, external bodies or others in the collection of materials and data, the design and construction of apparatus, the performance of experiments, the analysis of data, and the preparation of the thesis (including editorial help). In addition, it is appropriate to recognize the supervision and advice given by the thesis supervisor(s) and members of TSC.

### **B.1.5 Abstract**

Every copy of the thesis must have an English abstract, being a concise summary of the thesis, in 300 words or less.

### **B.1.6 Bibliography**

The list of sources and references used should be presented in a standard format appropriate to the discipline; formatting should be consistent throughout.

**Sample pages** of both MPhil and PhD theses are provided here (MPhil / PhD), with specific instructions for formatting page content (centering, spacing, etc.).

## **B.2 Language, Style and Format**

### **B.2.1 Language**

Theses should be written in English.

Students in the School of Humanities and Social Science who are pursuing research work in the areas of Chinese Studies, and who can demonstrate a need to use Chinese to write their theses should seek prior approval from the School via their thesis supervisor and the divisional head.

If approval is granted, students are also required to produce a translation of the title page, authorization page, signature page, table of contents and the abstract in English.

### **B.2.2    Pagination**

1. All pages, starting with the Title page should be numbered.
2. All page numbers should be centered, at the bottom of each page.
3. Page numbers of materials preceding the body of the text should be in small Roman numerals.
4. Page numbers of the text, beginning with the first page of the first chapter and continuing through the bibliography, including any pages with tables, maps, figures, photographs, etc., and any subsequent appendices, should be in Arabic numerals.
5. Start a new page after each chapter or section but not after a sub-section.

*Note: That means the Title page will be page i; the first page of the first chapter will be page 1.*

### **B.2.3    Format**

1. A conventional font, size 12-point, 10 to 12 characters per inch must be used.
2. One-and-a-half line spacing should be used throughout the thesis, except for abstracts, indented quotations or footnotes where single line spacing may be used.
3. All margins—top, bottom, sides—should be consistently 25mm (or no more than 30mm) in width. The same margin should be used throughout a thesis. Exceptionally, margins of a different size may be used when the nature of the thesis requires it.

### **B.2.4    Footnotes**

1. Footnotes may be placed at the bottom of the page, at the end of each chapter or after the end of the thesis body.
2. Like references, footnotes should be presented in a standard format appropriate to the discipline.
3. Both the position and format of footnotes should be consistent throughout the thesis.

### **B.2.5 Appendices**

The format of each appended item should be consistent with the nature of that item, whether text, diagram, figure, etc., and should follow the guidelines for that item as listed here.

### **B.2.6 Figures, Tables and Illustrations**

Figures, tables, graphs, etc., should be positioned according to the scientific publication conventions of the discipline, e.g., interspersed in text or collected at the end of chapters. Charts, graphs, maps, and tables that are larger than a standard page should be provided as appendices.

### **B.2.7 Photographs/Images**

1. High contrast photos should be used because they reproduce well. Photographs with a glossy finish and those with dark backgrounds should be avoided.
2. Images should be dense enough to provide 300 ppi for printing and 72 dpi for viewing.

### **B.2.8 Additional Materials**

Raw files, datasets, media files, and high resolution photographs/images of any format can be included.

*Note: Students should get approval from their department head before deviating from any of the above requirements concerning paper size, font, margins, etc.*