

# Extracting the Common Structure of Compounds to Induce Plant Immunity Activation using ILP

Atsushi matsumoto,<sup>1,2</sup> Katsutoshi Kanamori,<sup>1</sup> Kazuyuki Kuchitsu,<sup>3</sup>  
and Hayato Ohwada<sup>1</sup>

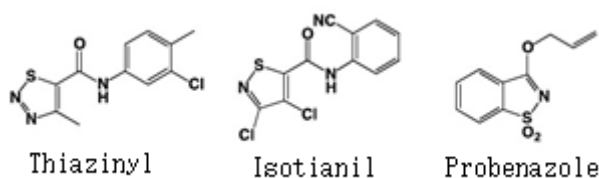
1. Department of Industrial Administration, Faculty of Science and Technology,  
Tokyo University of Science  
2. 7415617@ed.tus.ac.jp
3. Department of Applied Biological Science, Faculty of Science and Technology,  
Tokyo University of Science

**Abstract.** While recent studies have referred to plant immunity activators, it is difficult to find a compound to use for the immunity activation of plants. In this study, we seek to determine compounds that enable plant immunity activity using ILP. With the proposed method, it is possible to predict compounds that induce plant immunity activity, based on the structural features of the compounds. The predicted structure rule also includes structures of known plant immunity activators. However, further investigation is needed regarding the relationship between plant immunity and structure rules.

**Keywords:** ILP, Machine learning, Plant immunity activation, Virtual screening

## 1. INTRODUCTION

Decreased production of agricultural crops due to pathogenic bacteria and pests is a serious problem that has not yet been solved. To address this problem, grower have made a deal with fungicides and pesticides, however, it is difficult to act selectively on the target (e.g., pests and pathogens). There is a possibility that the cause of health damage in humans and destruction of biota. In addition, long-term use of the same drug may cause the emergence of resistant bacteria; thus, the effect of the drug gradually decreases. In recent years, plant immunity activators have attracted attention, based on the idea of increasing the immunity of the plant rather than directly killing pathogens and pests. However, only three types of plant immunity activator are currently marketed in Japan (Fig. 1). In addition, the mechanism of plant-immunity activation is still largely unknown [1].



**Fig. 1.** Known plant-immunity activators

The development of plant immunity activators has been slow, due to the time required and the high cost of screening candidate compounds. Cause of this problem is the kind of candidate compounds is enormous and each of the compounds were reacted to the cells to confirm the effect of immunity activation.

Therefore, virtual screening using machine learning has attracted attention in the field of drug discovery. In this study, we predict compounds that induce plant-immunity activation using ILP to study compound structures. ILP can be used to determine relationship patterns between data; therefore, it is suitable to represent the structure of compounds. Additionally, we obtained the structure of the predicted compound as a rule, which is one of the excellent points of ILP. A recent study that was conducted to predict the structure of compounds using ILP exhibited high performance [2]. In those cases, the target of compound bonds was known. However, in the present study, the target of compound bonds is not known.

## 2. PLANT IMMUNITY

Plant immunity is a defense system to protect plants from various enemies. A plant-immunity activator is a drug that activates plant immunity. The Kuchitu group constructed a screening system to find a candidate using the amount of ROS (reactive oxygen species) generation as an index [3]. Experiment results indicated that if the ROS value is high, the compound is likely to be a plant-immunity activator.

## 3. DATASET

In the present study, the datasets are experiment data about the plant immunity activator in *Arabidopsis thaliana*, compiled by the Kuchitu group. This dataset includes 10000 compounds. Positive examples are 271 high-ROS compounds, and negative examples are the other 9729 compounds. However, negative examples were reduced to 813 compounds by random sampling for two reasons. First, imbalanced data deteriorates learning accuracy. Second, if there are many compounds, calculation takes a long time. Therefore, 1084 compounds were used in this study (Table 1).

**Table 1.** Number of compounds used

positive	negative	total
271	813	1084

## 4. METHOD

### 4.1 ILP Approach

With the ILP approach, structural features and some numerical features of the compound were used as background knowledge. In this study, we used GKS [4], which is an ILP system. We defined seven predicates to represent the features of the compounds:

- `atom (compound_name, atom_id, element)`  
Types of atoms present in the compound
- `bond (compound_name, atom_id, atom_id, bondtype)`  
Bonding state between atoms and bond type in the compound
- `Num_AromaticRings (compound_name, Num_AromaticRing)`  
The number of aromatic rings in the compound
- `Num_Rings (compound_name, Num_Ring)`  
The number of rings in the compound
- `LogP98 (compound_name, value)`  
Lipid solubility of the compound
- `LogD (compound_name, value)`  
Indication of a change in lipid solubility by a change in Ph value
- `ring (compound_name, ring_id, atom_id, ringsize, ringtype)`  
Type of ring structure that is composed of each atom. It can represent the connection of the ring structure and other structures by using this predicate.

By selecting several predicates as background knowledge, we can obtain the structure of the compound as a learning result (Table 2). Background knowledge is a set of atomic formulas of each predicate. *Atom* and *bond* are always necessary. The reason why selecting *LogP98* and *LogD* is result of importance calculation using the average Gini coefficient.

**Table 2.** Predicates selected for background knowledge

Setting name	Predicate
ILP1	atom,bond
ILP2	atom,bond,Num_AromaticRings
ILP3	atom,bond,Num_AromaticRings,Num_rings
ILP4	atom,bond,ALogP98
ILP5	atom,bond,Num_AromaticRings,Num_rings,ALogP98,LogD
ILP6	atom,bond,Num_AromaticRings,Num_rings,LogD
ILP7	atom,bond,LogD,ring
ILP8	atom,bond,ring

#### 4.2 SVM Approach

We also tried SVM [5] for comparison with ILP, using 77 features for learning (Table 3).

**Table 3.** Attributes used for SVM

Types of features	The number of features
Related to structure	39
Related to ALogP	6
Related to size or weight	14
Related to energy	12
Other	6
Total	77

Cost parameters and gamma parameters were determined using a grid search for 20 split from 0.0001 to 10,000. The kernel used was RBF.

#### 4.3 Evaluation

Ten-fold cross-validation was used in both approaches.

## 5. RESULTS

Table 4 presents the ILP results.

**Table 4.** ILP results

Setting name	tp	fn	tn	fp	Accuracy	Precision	Recall	F value
ILP1	92	179	699	114	0.73	0.447	0.339	0.386
ILP2	116	155	644	169	0.701	0.407	0.428	0.417
ILP3	127	144	605	208	0.675	0.379	0.469	0.419
ILP4	88	183	712	101	0.738	0.466	0.325	0.383
ILP5	131	140	572	241	0.649	0.352	0.483	0.407
ILP6	139	132	568	245	0.652	0.362	0.513	0.424
ILP7	165	106	523	290	0.635	0.363	0.609	0.455
ILP8	165	106	542	271	0.652	0.378	0.609	0.467

Table 5 compares the best of SVM and the best of ILP

**Table 5.** Comparison of the best of SVM and the best of ILP

Approach	tp	fn	tn	fp	Accuracy	Precision	Recall	F value
SVM	123	148	703	110	0.762	0.528	0.454	0.488
ILP8	165	106	542	271	0.652	0.378	0.609	0.467

Table 6 lists the best rules obtained using ILP. A good rule has many positive examples and few negative examples.

**Table 6.** Rules for compound structure

Rule number	Interpretation	Positive	Negative
Rule1	Atom C has a single bond with the aromatic ring.	27	10
Rule2	There is an aromatic ring containing an atom S and atom C has a double bond with something.	20	8
Rule3	Two aromatic rings bond to each other and each aromatic ring have a single bond.	22	10
Rule4	An aromatic ring containing an atom N and An aromatic ring consisted of 5 atoms bond to each other	15	3
Rule5	An aromatic ring containing an atom S and another aromatic ring bond to each other	14	2

## 6. CONCLUSION

Although SVM F values slightly exceeded those of ILP, ILP tp values greatly exceeded those of SVM. For virtual screening, it is very important to reduce the positive example of misclassification. Results of this study indicate that structural features of the compounds are useful in predicting immunity activation.

Using the ring structure as background knowledge yielded better results than not using ring structure. Therefore, the ring structure is considered an important factor in plant immunity activation.

When analyzing rules using ILP, comparison of known plant immunity activators indicated that Rule 2 was true for all three compounds. For rule showing a structure that is different from the known plant immunity activator, there is a need for further investigation.

In this study, it was possible to predict the partial structure that exists in all compounds of known plant-immunity activators. In addition, the rule that is unknown the relationship between immunity activity has been predicted. In order to improve prediction accuracy, it is essential to improve background knowledge in the future.

## References

1. Yoshiteru Noutoshi, Masateru Okazaki, Tatsuya Kida, Yuta Nishina, Yoshihiko Morishita, Takumi Ogawa, Hideyuki Suzuki, Daisuke Shibata, Yusuke Jikumaru, Atsushi Hamada, Yuji Kamiya, Ken Shirasu, Novel Plant Immune-Priming Compounds Identified via High-Throughput Chemical Screening Target Salicylic Acid Glucosyltransferases in Arabidopsis. *The Plant Cell*, vol.24:3795-3804, 2012
2. Jose C A Santos, Houssam Nassif, David Page, Stephen H Muggleton, Michael J E Sternberg, Automated identification of protein-ligand interaction features using Inductive Logic Programming:a hexose binding case study. Santos et al. *BMC Bioinformatics* 2012, 13:162, 2012
3. T Higashi, T Kurusu, S Hasegawa, K Kuchitsu, Dynamic intracellular reorganization of cytoskeletons and the vacuole in defense responses and hypersensitive cell death in plants. *Journal of Plant Research*, Volume 124, Issue 3, pp315-324, 2011
4. Hayato Ohwada, Hiroyuki Nishiyama, Fumio Mizoguchi, Concurrent execution of optimal hypothesis search for inverse entailment. *Lecture Notes in Artificial Intelligence*, Springer-Verlag, No.1866, Vol.4, pp.165-173, 2000
5. V.Vapnik, *The Nature of Statistical learning Theory*. Springer-Verlag, NY, USA, 1995