

CS 181 Machine Learning

Practical 1 Report, Team *Advanced Representation*

Jing Wen¹, (Jeremiah) Zhe Liu¹, and (Vivian) Wenwan Yang²

¹Department of Biostatistics, Harvard School of Public Health

²Department of Computer Science, Harvard School of Engineering and Applied Sciences

March 6, 2015

1 Exploratory Analysis

The training data set features the malware class (14 discrete categories) and XML behavior log regarding 3086 executables.

The empirical distribution of malware classes is displayed in Figure 1. As shown, majority of malwares are of relatively low frequency ($< 2\%$), except *Swizzor* (542, 17.56%), *VB* (376, 12.18%), and *Agent* (114, 3.69%).

Agent 114	AutoRun 50	FraudLoad 37	FraudPack 32	Hupigon 41	Krap 39	Lipler 53	
Magania 41	None 1609	Poison 21	Swizzor 542	Tdss 32	VB 376	Virut 59	Zbot 40

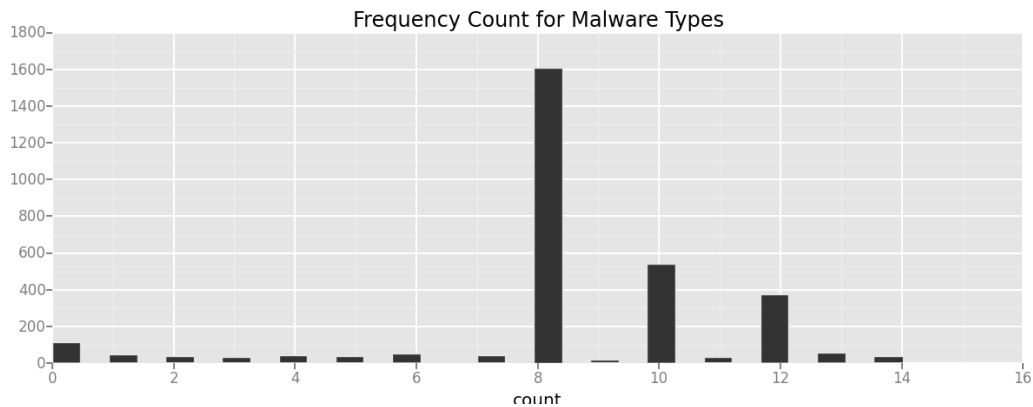


Figure 1: Empirical Count of Malware Types

The Input space is comprised XML behavior logs of the executables, which are generated from sand-box type environment where all system calls (as well as associated parameters and outcome status) of the executable are intercepted and reported in a structured XML file. Due to the complicated nature of the input space (structured strings), in the next section we discuss possible representation of such space and how to embed it into a euclidean/categorical feature space.

2 Embedding and Feature Extration

2.1 Matrix Representaion of Extraction Target

```
<load_image \
  filename="c:\342c547b28e9517f6fcf6c703933c0d9.EX" \
  successful="1" address="&#x24;400000" \
  end_address="&#x24;414000" size="81920" \
  filename_hash="hash_error"/>
```

Each XML log is consisted of a ordered set of structured report of system calls. As shown above, each call is consisted of a **operation type** (e.g. `load_image`), and **argument blocks** (e.g. `filename`, `successful`, `end_address`, `filename_hash`).

We currently consider information contained in such XML file as a ordered set of python dictionaries. To reasonable represent such information in numeric form, Trinius et al (2009) argued that the parameters in each call can be stratified into multiple levels, with the supposedly most informative information (e.g. operation type) to be contained level 1, some less varying argument blocks (e.g. the extension of loaded dll in `load_dll`) to be contained in level 2, and the more varying, call-specific (thus less informative) argument blocks to be contained in level 2 and beyond. This is the so-called Malware Instruction Set (MIST) representation.

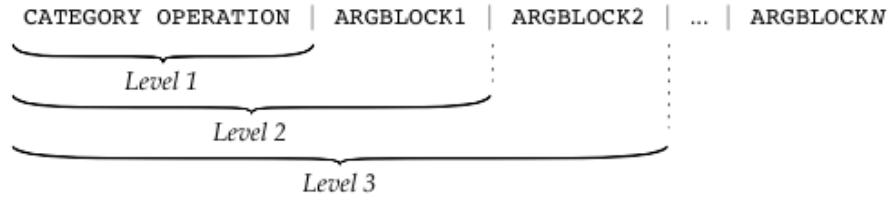


Figure 2: Malware Instruction Set (MIST) representation of XML records

In order to represent both the information contained in each call and also the order of the system call, we not only represent call-specific features, but also rep

2.2 Feature Extraction and Prunning

Order: Eigenvalue/Eigenvector

Call Counts: Row Sum

3 Model Choice and Justifications

3.1 Technical Detail of Random Forest

4 Model Choice and Justifications

Based on observations from previous section, we realize that current learning task possesses below characteristics:

1. Large Sample Size
2. High dimensional, binary Input Space
3. Complicated and unknown Feature-Outcome relation
4. Prediction as the Learning Goal

Based on above characteristics, we considered four classes of popular machine learning strategies as our potential model: regularized linear model (Ridge/LASSO), kernel methods (SVM/Gaussian Process), Bayesian methods, and random forest.

We decided against using **regularized linear models** due to the biased nature of the estimators. For example, recall the form of the likelihood function and solution to Ridge regression:

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{T} - \mathbf{X}\mathbf{w}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|_p^2$$
$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}$$

In order to specify a reasonable linear model, we have to include all possible n-level interactions between the 31 non-trivial features, which will necessarily result in a non-trivial λ during the process of penalization. However, if the target values truly follow linear form, i.e. $\mathbf{t}_{\text{True}} = \mathbf{X}\mathbf{w}_{\text{True}}$, we have:

$$\hat{\mathbf{t}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T * (\mathbf{X}\mathbf{w}_{\text{True}}) \neq \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T * (\mathbf{X}\mathbf{w}_{\text{True}}) = \mathbf{t}_{\text{True}}$$

A simple eigen decomposition will illustrate that such bias increase with the magnitude of λ . We hence decide not to use regularized linear models since the interest lies in unbiased prediction of a potentially complex relationship. In practise, Ridge estimator gives Average Mean square Error (AME) around 0.29, offering inferior performance compared to the other algorithms.

We also decided against using **Kernel-based methods** due to the large sample size. This is because instead of working with the design matrix, Kernel methods working with a $n \times n$ Kernel matrix of the form $\mathbf{K} = \Phi\Phi^T$, where Φ is a $n \times p$ kernel-transformed matrix of input features. Applying kernel matrix in current scenario implies working with a $1,000,000 \times 1,000,000$ potentially non-sparse matrix, whose computation cost is prohibitive for obvious reasons.

We also chose not to use **Bayesian Methods** due to our lack of prior knowledge with respect to the effect of each features. As a result, Bayesian regression doesn't seem to provide any benefit over regular models.

As a result, we are left to choose **Random Forest** as our algorithm of choice. Introduced by Breiman and Cutler in 2001, random forest is essentially a ensemble method for regression trees. It offers a relatively scalable algorithm which is robust against noisy observations / covariates and powerful in capturing complex interaction between features. This method is traditionally considered having superb performance in terms of prediction accuracy.

4.1 Technical Detail of Random Forest

In this practical, we adopted the implementation `RandomForestRegressor` from `scikit-learn 0.15.1`, which offers a paralleled implementation of the original algorithm described in Breiman (2001). We describe the pseudo code for this algorithm in **Appendix**, and discuss in this section several important parameters in the algorithm:

B (Number of Trees in Forest).

One of the important features of random forest is it does not overfit, i.e. increasing the number of trees in forest does not overfit the data. Indeed, random forest estimate approximate the expectation of "real tree" conditional on the space of bootstrap samples $\hat{f}_{rf}(x) = E_{\Theta(Z)} T(x|\Theta(Z)) = \lim_{B \rightarrow \infty} \hat{f}(x)_{rf}^B$ (Hastie et al, 2009). Increase B will lead to a less biased and less noisy approximation of such value.

m (Number of Sampled Features) & d_{max} (Maximum Tree Depth).

Although increase **B** will only bring estimator toward a constant limit, this limit itself, however, may overfit depending on **d_{max}** and **m**. This is because intuitively, too deep a tree indicates a overly rich model, and too many considered feature increase correlation between generated trees. Both may incur unnecessary variance. It is thus crucial to perform parameter selection for **m** and **d_{max}** during model fitting.

5 Parameter Selection

5.1 Important Features

As discussed in Section 1 (Exploratory Analysis). Only 31 features were considered important for the purpose of prediction. To validate this point, we first ran a naive random forest with default setting (**m** = **p**, **d_{max}** = ∞) and 100 trees, and consider the variable importance metric produced by this fit. Since the variable important for a specific feature is constructed by considering in every tree the change in out-of-bag prediction error by setting the "effect" of this feature to 0, it is considered a cross-validated metric of the contribution of each features in terms of prediction. As a result, 227 features have variable importance $\leq 1e-5$ and were thus removed from the subsequent analysis. The 10-fold CV AME of the naive model is 0.272452, and we seek to improve our random forest model based upon this baseline.

5.2 Tunning Parameters

In the next step, we consider the effect of **m** and **d_{max}** in our prediction performance through 10-fold cross validation. With AME as the outcome metric, we seek to search over a grid of (**m**, **d_{max}**) since the convexity of the current problem is unclear. For regression problems, Breiman recommended setting **m** = $\frac{p}{3}$ and node size equal to 5, which equals to **m** = $28/3 \approx 10$ and **d_{max}** = $\log_2(1000000/5) \approx 17.6$. In practice, however, larger **m** usually works better for more complex problems. We hence decide to search over the grid $[\mathbf{m} \in \{14, 16, \dots, 28\}] \times [\mathbf{d}_{max} \in \{16, \dots, 24\}]$ so that it covers a reasonable range of the theoretically sound values. The AME surface over candidate parameter values are visualized in Figure ?? and Table ?. As a conclusion, we selected (**m** = 19, **d_{max}** = 28). With the selected parameters, we ran a random forest with size **B** = 2000 in order to achieve closest approximation toward the theoretical limit $E_{\Theta(Z)} T(x|\Theta(Z))$ within reasonable machine computing time (5 hours). The final 10-fold CV AME within training dataset is calculated to be 0.271845.

6 Model Assessment and Discussion

The observed gap value v.s. prediction bias (absolute value of residual) is visualized in Figure ?? (right). As shown, the prediction bias increases linearly when gap values are far from its mean, which is expected given the limited number of category combinations (4546 unique combinations compared to 1,000,000 observations) of extracted features, which only allowed the random forest to divide the input space into as many as 4546 subspaces and predict using local mean of each subspace, which is why clusters of linearly increasing prediction bias are observed, indicating underfit due to limitation of the input space.

Given the pre-extracted feature provided by this dataset, authors believe that our current model should offer best possible performance in terms of prediction. However, prediction performance may be further improved by extracting more energy-relevant, and preferably continuous features from the molecular structure. (e.g. using **RDKit** to extract from the SMILES string), and explore feature-outcome relationship using again random forest to achieve a more satisfactory prediction result. Unfortunately, limited by the time and computing environment available to authors, above procedure was not carried out.

Reference

1. L Breiman, Random Forests, Machine Learning, 45(1), 5-32, 2001.
2. T Hastie. R Tibshirani. J Friedman, "Elements of Statistical Learning", Springer, 2009.
3. P Trinius et al, "A Malware Instruction Set for Behavior-Based Analysis", Technical Report, University of Mannheim, Germany, 2009

Appendix: Algorithm for Random Forest

Input: $Z = (t_{N \times 1}, X_{N \times p})$

Parameter:

- B : The number of trees in the forest
- m : The number of features to sample when splitting nodes.
- d_{\max} : The maximal allowed depth of each tree.
- n_{\min} : The minimal allowed size the of the splitted node.

Algorithm:

1. For $b = 1$ to B
 - (a) Draw a bootstrap sample Z^* of Z
 - (b) (Grow a regression tree T_b using Z^*)
Untile maximal depth of tree (d_{\max})/minimal size of node (n_{\min}) is reached, for every node in current tree:
 - i. Randomly sample m features from the p features
 - ii. Find optimal split (in the sense of minimizing RMSE) of current node with respect to select features.
 - iii. Split the node into two daughter nodes.
2. Output the ensemble of trees $\{T_b\} \ b \in \{1, \dots, B\}$.
3. Predict for a new point \mathbf{x} as: $\hat{f}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x})$.