

CS 181 – PRACTICAL 3

CASEY GRUN, SAM KIM, RHED SHI

1. WARMUP

We implemented multiclass logistic regression and generative classifiers with Bayesian class-conditional densities on fruit data consisting of 59 data points sorted into apples, oranges, and lemons. The number of classes is $K = 3$. The data points reflect the height and width of each fruit, giving us two parameters \mathbf{x}_1 and \mathbf{x}_2 . Our basis function is given by

$$\Phi(\mathbf{x}) = w_0 + w_1 \cdot \mathbf{x}_1 + w_2 \cdot \mathbf{x}_2 = \mathbf{w}^T \mathbf{x}$$

1.1. Logistic Regression. We plotted the coloring of the graph using the maximum $p(C_k)$ for each point in the coordinate space. Admittedly, our algorithm fails to capture the lemons, perhaps due to the algorithm not running long enough. We randomized the guesses for our \mathbf{w} vector and ran the optimization using the inverse Hessian and gradient of the error function given by

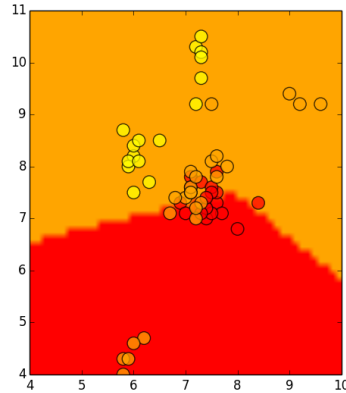


FIGURE 1. Logistic regression with decision boundaries.

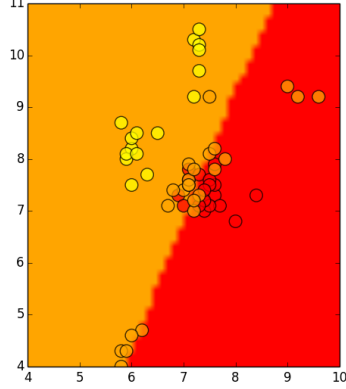


FIGURE 2. Generative classifier with decision boundaries.

$$\begin{aligned}\nabla E(\mathbf{w}) &= \Phi^T(\mathbf{y} - \mathbf{t}) \\ \mathbf{H} &= \Phi^T \mathbf{R} \Phi\end{aligned}$$

1.2. Generative Classifier with Bayesian Class-Conditional Densities. We defined the weights of our generative classifier method with Bayesian class-conditional densities as

$$\begin{aligned}\mathbf{w}_k &= \Sigma_k^{-1} \mu_k \\ w_{k0} &= -\frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k + \ln(p(C_k))\end{aligned}$$

where we have

$$\begin{aligned}\mu_k &= \frac{1}{N_k} \sum_{n=1}^N t_n \mathbf{x}_n \\ \Sigma_k &= \frac{1}{N_k} \sum_{n \in C_k} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T\end{aligned}$$

Our results are shown here in the figure. Although there only appears to be one decision boundary between the apples and oranges, this boundary shows curved features not seen in the logistic regression because of the individual covariances in our weight calculations.

2. CLASSIFICATION

Our challenge was to classify a set of programs, based on traces of their system calls, as either a type of malware or a non-malicious program. Given an $N \times D$ matrix \mathbf{X} representing the programs, along with a N -dimensional vector \vec{t} of training labels, produce a function

$$y : \mathbf{X}' \mapsto \vec{t}'$$

. That is, a function which could produce labels \vec{t}' for some matrix \mathbf{X}' . There were two parts to this challenge: determining the *feature functions* to generate rows of the matrix \mathbf{X} based on the system calls for each malware program, and determining what classification algorithm to use to generate the function y .

We developed a cross-validation procedure to test our predictions on folds of the training data before submitting our predictions to Kaggle on the test data; we used 10 folds. A persistent problem we observed was that predicted scores in cross-validation were 0.1–0.2 higher than the scores we received on the public leaderboard on Kaggle. We suspect this was because of the heterogeneous training data

2.1. Feature Functions. We evaluated a number of feature functions:

System call counts: Our first feature simply summed the number of system calls across all threads in all processes for a particular program. This method achieved 41.6% on the public leaderboard.

n -grams of system calls: Our most substantial feature examined counts of sequences of system calls. Each unique sequence of n system calls was represented as a feature that counted how many times this sequence appeared in any given malware. We used cross-validation to try and optimize the choice of n —the length of these sequences. Too small an n would result in uninformative sequences, while too large an n would not generalize well. We found the optimum n to be approximately 8.

DLLs loaded: We also created features counting how many times each unique DLL was loaded with the `load_dll` system call. This feature proved to be reasonably informative, and achieved our best score on Kaggle (in combination with the 8-grams of system calls).

Registry keys: Finally, we employed a similar method to count the number of times each unique registry key was read by a program. This feature was less informative than the DLLs or 8-grams, so it was not included in our final predictions.

Additionally, we attempted classification on a reduced-dimensionality feature set, generated by performing a truncated singular value decomposition (TSVD) on the full feature set. This method did not perform particularly well, so it was not pursued.

Our final/best score was achieved by combining the DLLs, 8-grams, and system call counts features.

2.2. Classification Methods. We evaluated a number of classification learning algorithms, as implemented by the scikit-learn Python library. These include:

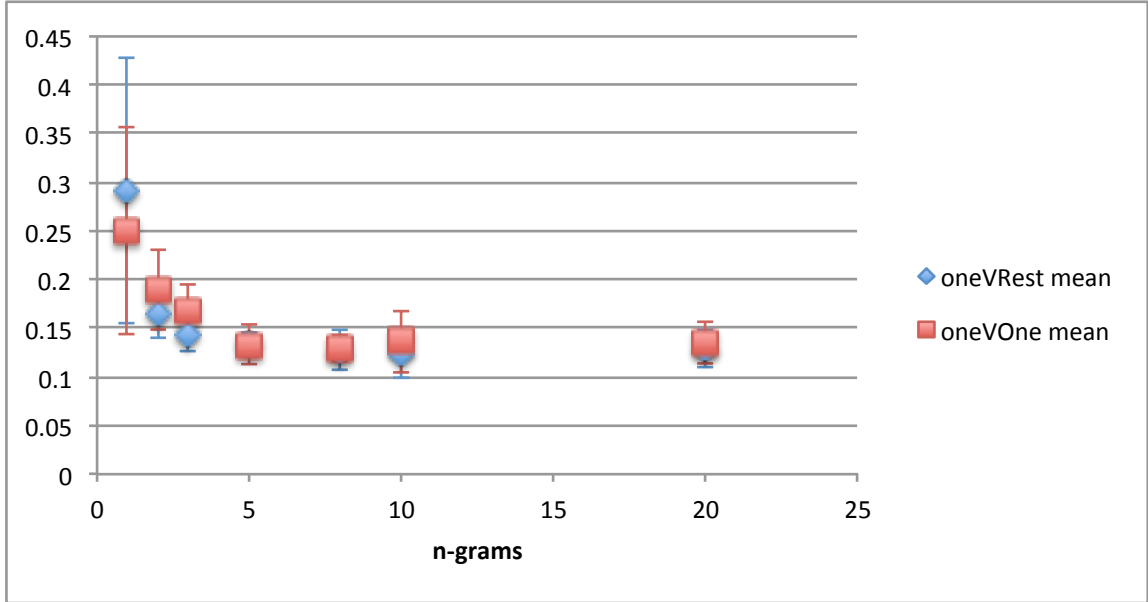


FIGURE 3. Tuning the length n of sequences of system calls. Y-axis is error.

One-vs-rest: As discussed in the textbook, this algorithm fits one classifier per class to classify if a point is in a class C_k or not. We used linear discriminant functions.

One-Vs-One: This algorithm introduces a discriminant function for every pair of classes, where each point is classified according to a majority of decisions among the discriminant functions. We used linear discriminant functions.

Logistic regression: This algorithm minimizes an error based on the training set to optimize the parameters for the likelihood function.

Decision trees: This algorithm goes down a tree of binary nodes in which the inputs are used to determine which path to go down, which ultimately assign probabilities to the target variable.

Forests of randomized trees: Related to the decision trees, this algorithm uses a set of trees, where each tree is built using random decisions. Whereas a decision tree will use the best split among all the features for each node, the random tree will pick the best split from a random subset of features. The trees are built independently and then their results are averaged to make predictions. This results in the bias increasing (due to not picking the best split) but the variance decreasing (due to averaging) which offsets the increase in bias. In our testing, we used a forest of 10 trees, and a maximum number of features in the random subset of $\sqrt{\text{num_features}}$.

Forests of extremely randomized trees: This is similar to the forests of randomized trees algorithm, except that splits are computed slightly differently in building the random trees. In addition to using a random subset of features for each node, instead of picking the best feature, thresholds are drawn randomly for each feature

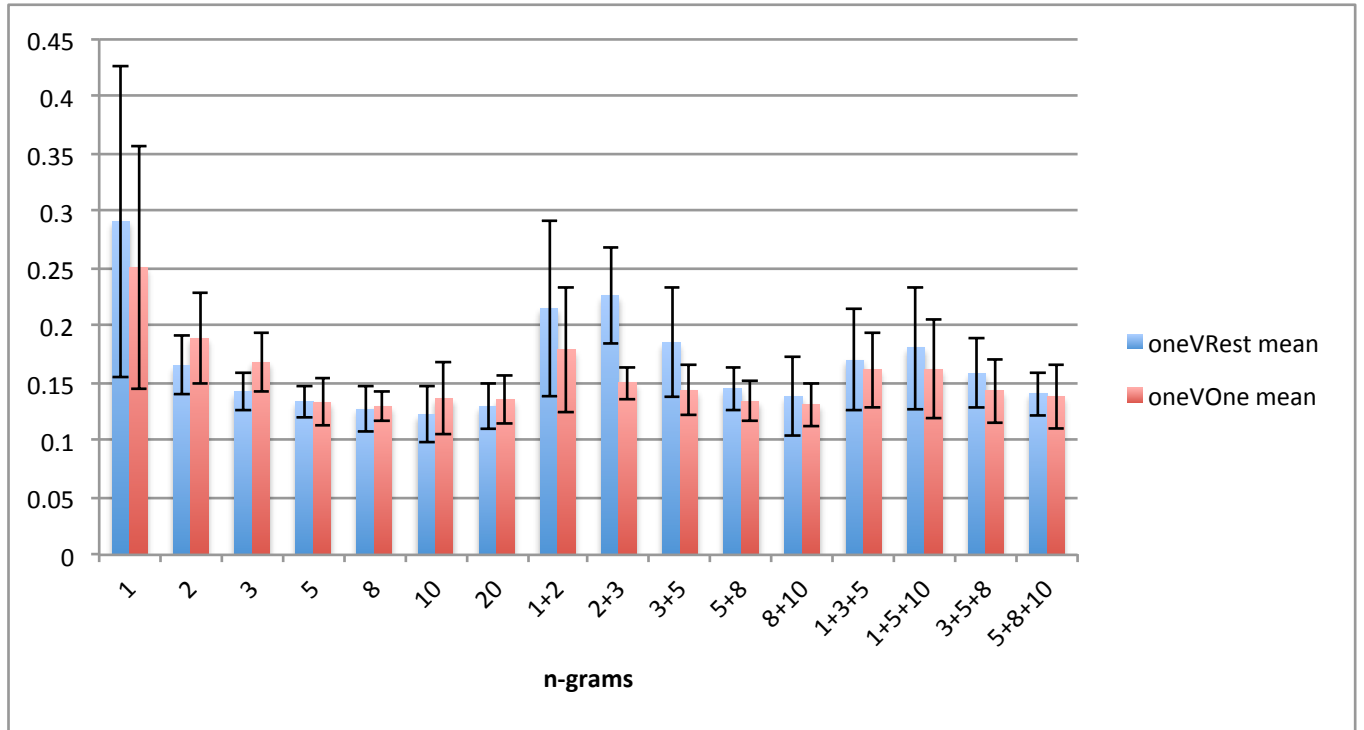


FIGURE 4. Considering combinations of n -grams of system calls. Y-axis is error.

and these thresholds are picked as the splitting rules. This results in an even higher bias but lower variance. In our testing, we used the same parameters as the forests of randomized trees.

We evaluated these learning algorithms based on the extracted features of system call counts and system call bigrams, and using cross-validation methods to withhold 10% of the data over 10 trials and averaging the percent errors over the predictions on the withheld data. Standard deviations over the 10 trials were calculated and the average running time for a single learning and prediction run are also noted. The results are summarized in the table below:

Classification algorithm	% error	Standard deviation	Running time (s)
One-vs-rest (linear)	12.7	2.4	7-8
One-vs-one (linear)	12.9	1.4	7-9
Logistic regression	11.5	1.5	600-700
Decision tree	10.9	1.3	50-60
Random trees	10.4	1.2	8-10
Extremely random trees	9.7	1.6	16-20

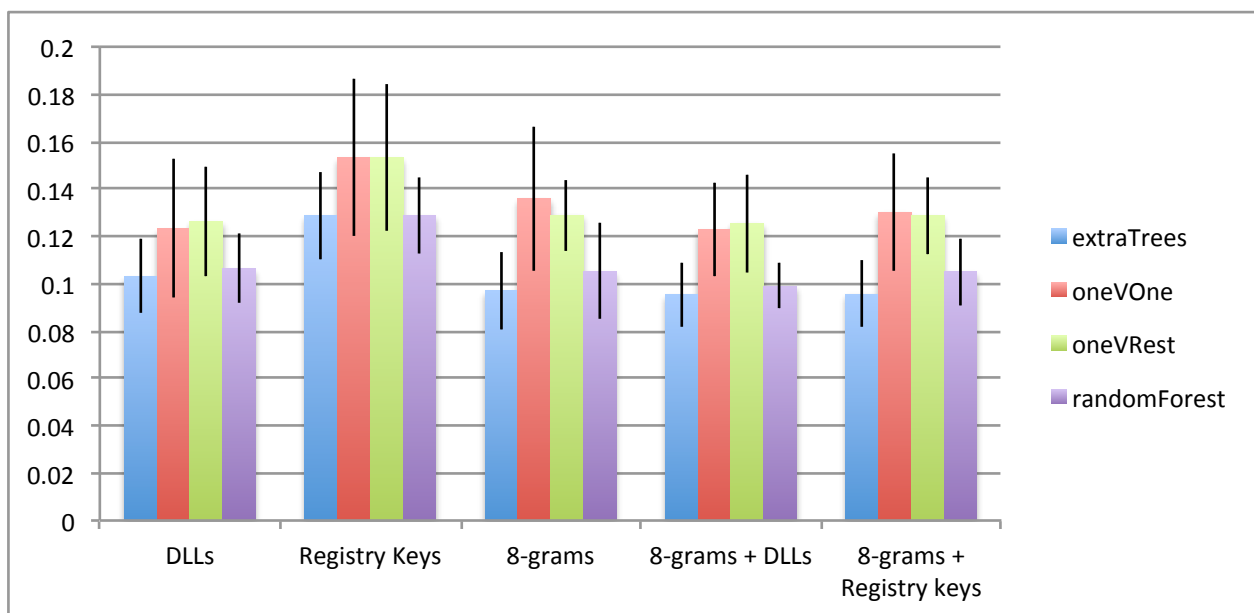


FIGURE 5. Comparison of best-performing features. Y-axis is error.

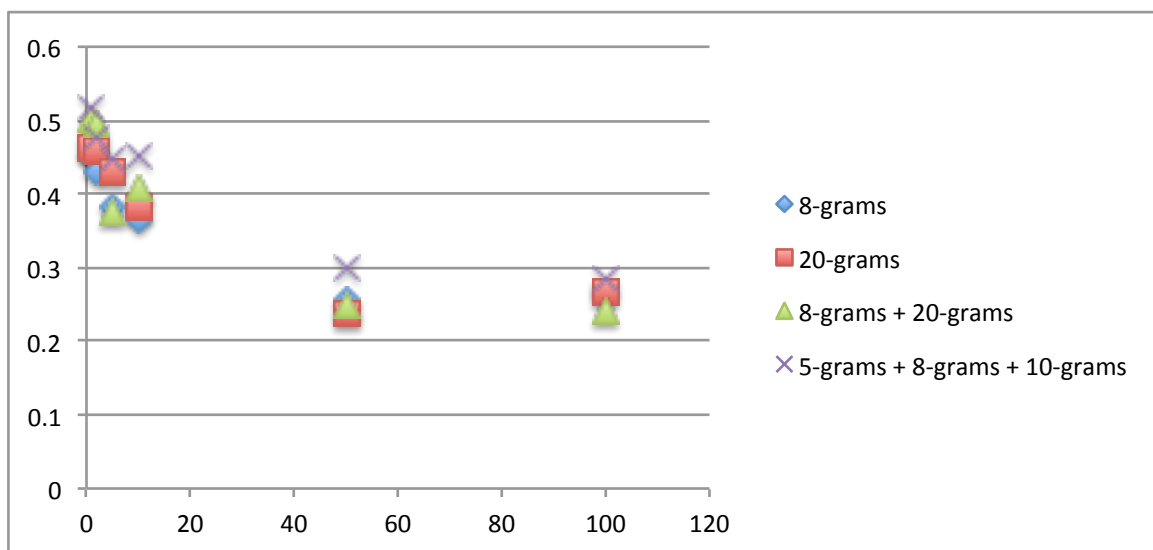


FIGURE 6. Comparison of a truncated singular value decomposition (TSVD) on various combinations of n -grams. The TSVD performed considerably worse than the full feature set in cross-validation, so it was not pursued.

Additionally, on Kaggle, the one-vs-rest, decision tree, and extremely random trees classification algorithms using the system call counts and system call bigrams reached accuracies of 68.5%, 79.1%, and 81.0%, respectively.

We suspect that decision trees and its respective variants did better than the other classification algorithms because decision trees are not limited to functional decision boundaries in the feature space. Granted, we only used linear discriminant functions for one-vs-rest and one-vs-one classification algorithms, and these performed surprisingly well. However, further investigation would require using different basis functions or boundaries. Decision boundaries are not limited to functional decision boundaries.

Our team name on Kaggle is "Cheesy Dramas"

REFERENCES

- [1] Scikit-learn: Machine Learning in Python, Pedregosa *et al.*, JMLR 12, pp. 2825-2830, 2011