# Binary Classification of Internet Traffic
## Identifying Distributed Denial of Service Attacks
## Chris Soyars

KENNESAW STATE UNIVERSITY
COLLEGE OF COMPUTING AND SOFTWARE ENGINEERING
School of Data Science and Analytics

Faculty Advisor: Professor Michael Frankel

## INTRODUCTION

- **Distributed Denial of Service (DDOS)** is a common attack method for Internet service disruption.
- **Multiple attack methods can result in DDOS**: broad detection and mitigation ability is a must to minimize vulnerability.
- **Attacks are often concurrent with legitimate traffic**: goal of mitigation is to filter malicious traffic while allowing benign traffic to pass without significant obstruction.
- Data set for building logistic regression contains 7,616,509 observations and 85 variables, including a Label variable for Benign vs DDOS traffic. 17% of observations are from DDOS attacks.

## METHODS

- **Examine data set** to gather comprehensive information on variable properties. Eliminate variables that cannot be used due to poor or lacking information. Identify and impute missing or erroneous data if possible.
    - ~50000 observations had values NA, infinity, or implausibly negative for certain variables due to lack of precision and accuracy in time-related information. Values were recalculated to restore information.
- **Cluster variables** to identify variables of greatest interest.
- **Discretize variables and calculate odds of benign vs ddos traffic for each bin** to extract additional information and trends from selected variables. Natural log of odds also calculated for each bin.
- **Eliminate non-significant variables** from selection for logistic model.
- **Create logistic regression model** using selected variables. Model is trained on 80% of the data set.
- **Refine logistic regression model** to eliminate redundancy and select for highly significant variables.

## RESULTS

- Logistic regression model with 17 selected variables has 99.9% concordant pair rate (C=0.999). Model predicts a lower probability of possible DDOS traffic for benign cases than actual DDOS attacks in most observations.
- Maximal KS statistic of .954
- Model can be simplified depending on needs and resources. C=.995 with as few as 4 variables.
    - Initial measurements only: C=.965 with 3 variables, with marginal gains for additional initial statistics.
- Sensitivity and specificity are maximized with probability threshold of 0.188.
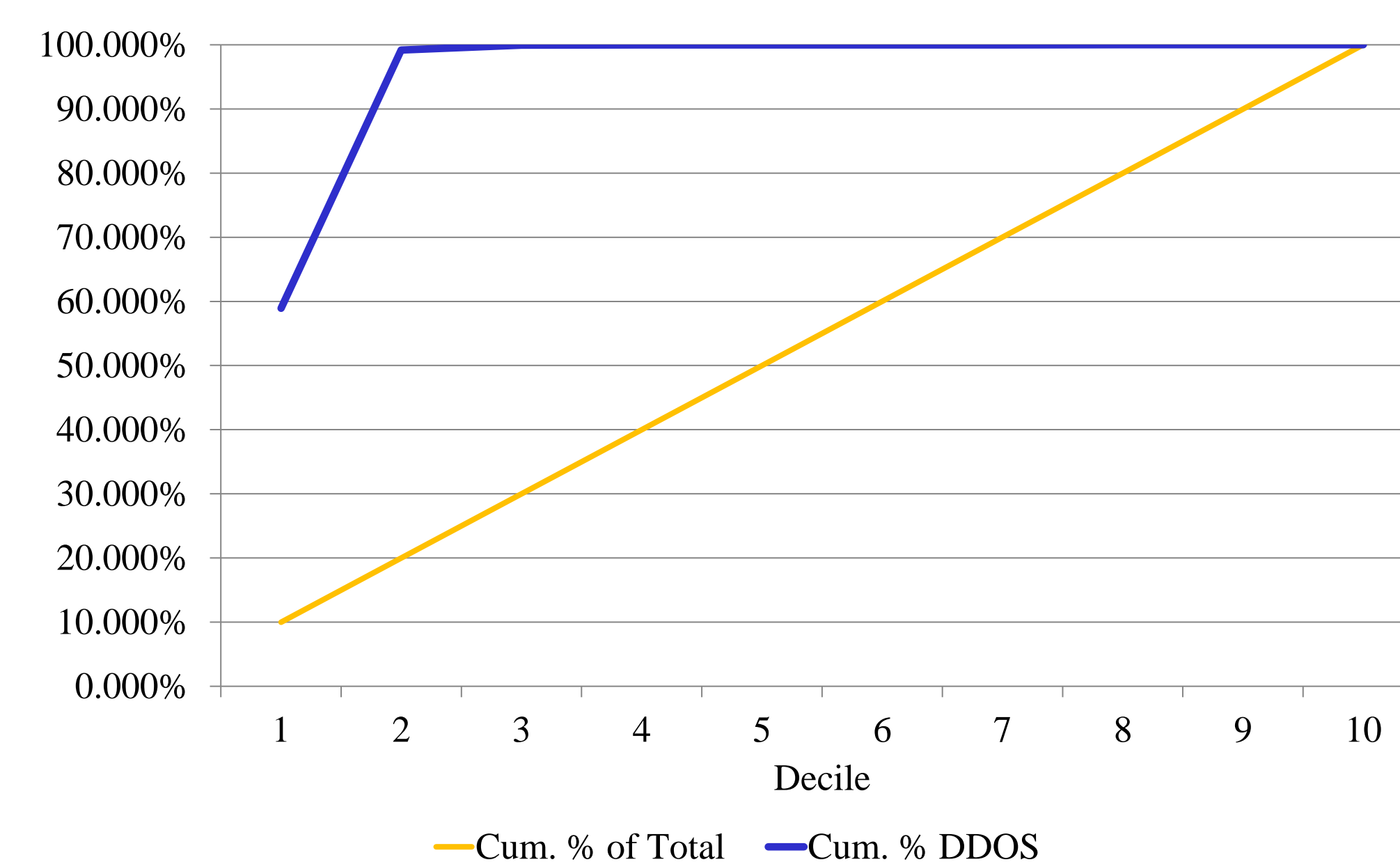    - 17 variable model correctly identifies 98.7% of benign traffic and 98.2% of DDOS traffic.

### Table 1. Concordance Statistics for Selected Models

| Number of Variables | Concordance |
|---|---|
| 17 | 0.999 |
| 12 | 0.998 |
| 4 | 0.995 |
| 3 (initial statistics) | 0.965 |

### Table 2. Confusion Matrix for Logistic Regression Model with 17 Variables

| Label Frequency Percent | Prediction | | |
|---|---|---|---|
| | Benign | DDOS | Total |
| Benign | 1248228 | 16168 | 1264396 |
| | 81.94 | 1.06 | 83.00 |
| DDOS | 4410 | 254495 | 258905 |
| | 0.29 | 16.71 | 17.00 |
| Total | 1252638 | 270663 | 1523301 |
| | 82.23 | 17.77 | 100.00 |

### Figure 1. KS Chart for Benign vs DDOS Classification



### Figure 2. ROC Curve for Logistic Regression Model



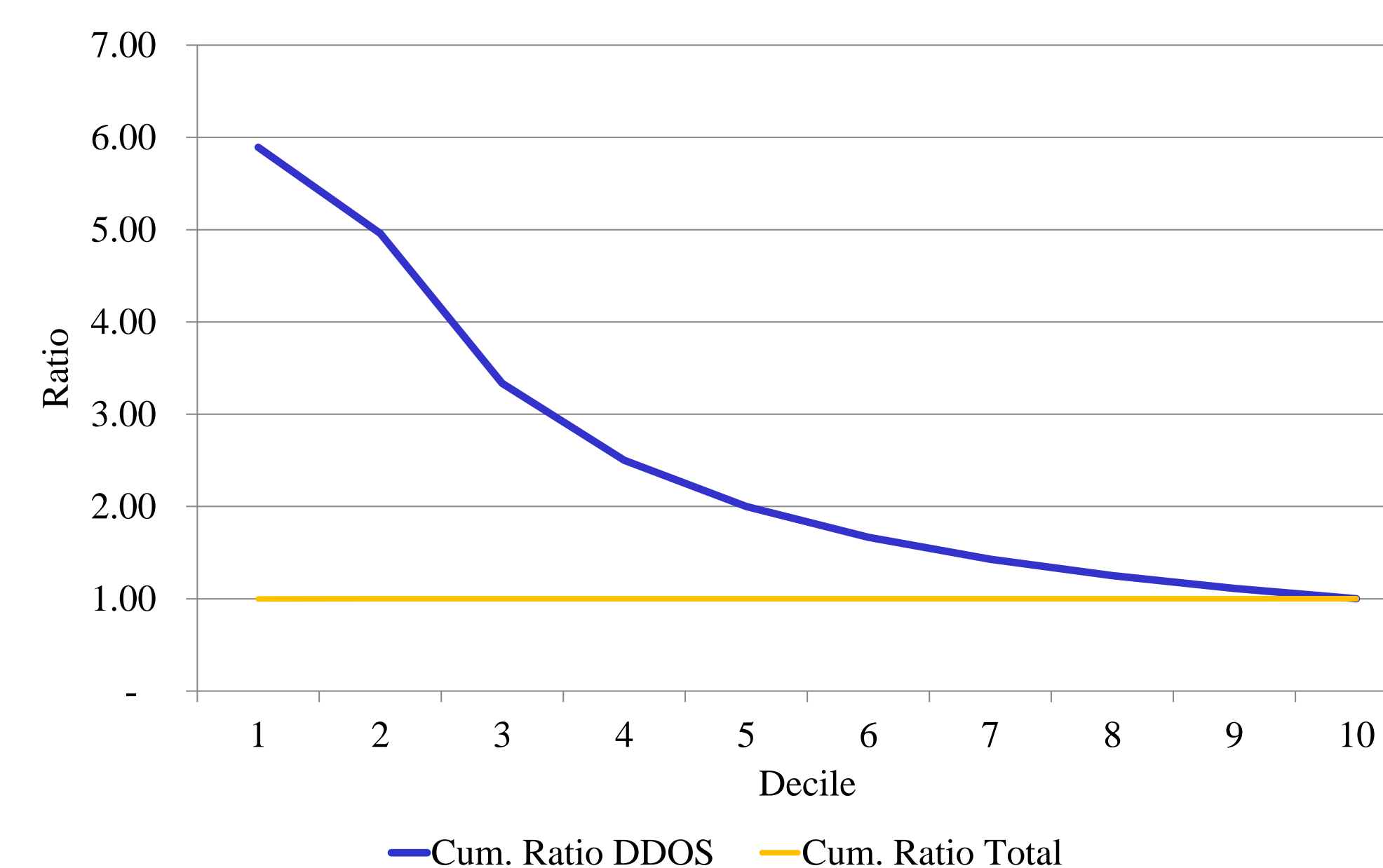### Figure 3. Gain Chart for Benign vs DDOS Classification



### Figure 4. Lift Chart for Benign vs DDOS Classification



## DISCUSSION

- **The selected logistic regression model has very high accuracy in classifying benign vs DDOS traffic.**
    - High separation between benign and DDOS traffic.
    - High success rate despite several variations in attack vectors and methods.
    - Effective as a first line of defense in filtering malicious traffic.
    - Content delivery and DDOS mitigation services have strong records in detecting attacks.
- **Predictive success in data set does not guarantee similar level of success moving forward.**
    - Traffic patterns, attack methods, and security practices can evolve rapidly.
    - Model likely to benefit from periodic re-evaluation.
    - Unknown success rate against novel attack vectors of the same general class.
- **Limitations & Improvements**
    - No analysis of trends by time or location (IP address)
    - Lack of precision on some variables
    - Possible errors from traffic logging software used

## SAS CODE

```
%IMPV5 (DSN=class.test, VARS=&varlist, EXCLUDE=Label, PCTREM=1,MSTD=);

PROC SQL;
    SELECT NAME INTO: VARNAME SEPARATED BY ' '
    FROM DICTIONARY.COLUMNS
    WHERE UPCASE(LIBNAME)="DDOS" AND
UPCASE(MEMNAME)="DDOS" AND NAME NOT IN("Label");
QUIT;
PROC VARCLUS DATA=import OUTTREE=tree MAXCLUSTERS=71;
    VAR &varname;
RUN;

PROC GLMSELECT DATA=ddos.disc2;
    MODEL label=&mvar / DETAILS=all SELECTION=lasso
STATS=all;
RUN;

PROC LOGISTIC DATA=train DESC OUTEST=betas
    OUTMODEL=scoringdata;
    MODEL label=&mvarsn /SELECTION=BACKWARD
    CTABLE pprob=(0.16 to 0.21 by 0.001)
    LACKFIT RISKLIMITS;
    OUTPUT OUT=output p=predicted;
    SCORE DATA=valid OUT=ddos.score;
RUN;
```

## ACKNOWLEDGEMENTS