# Linear Regression Analysis of Network Traffic Data

## Applied Machine Learning & Data Analytics

Dataset: Internet Firewall Logs
Professor: Dr. Samir Rustamov

Team Members:
Vusal Shirinbayli
Gabil Gurbanov

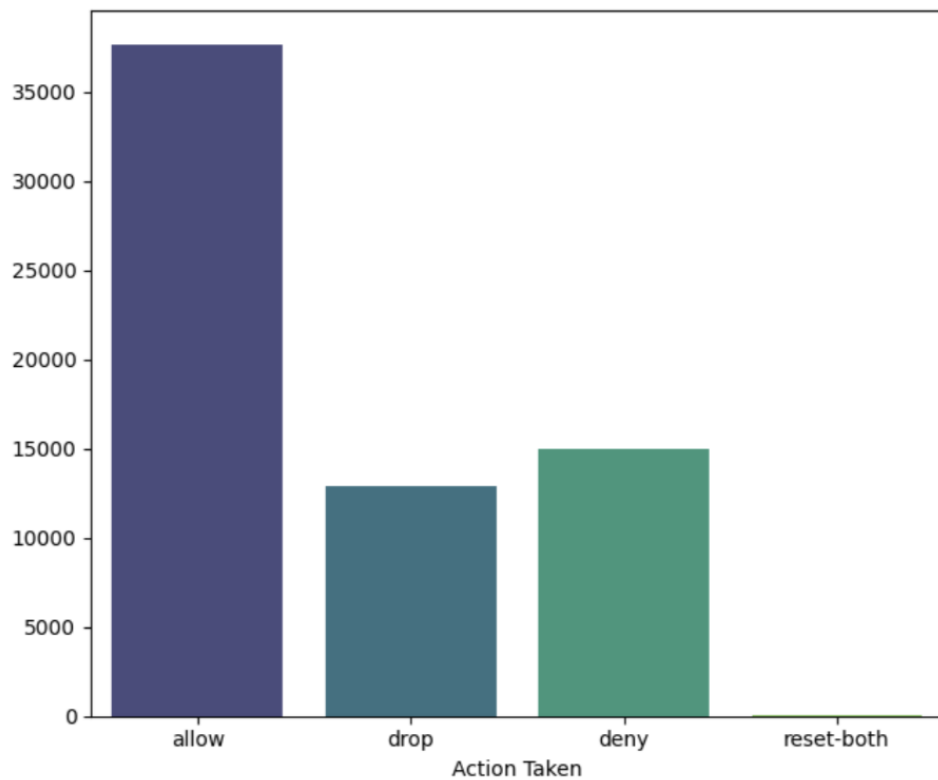# Dataset

Source: Kaggle Firewall Dataset
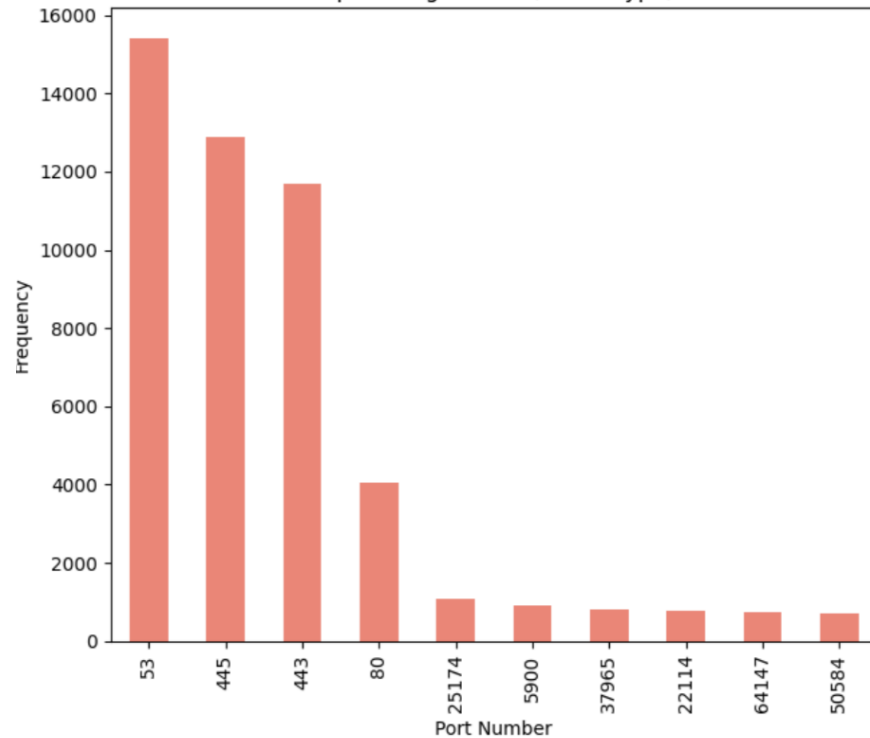
Observations: 65,532

Features: 12

Target: Bytes

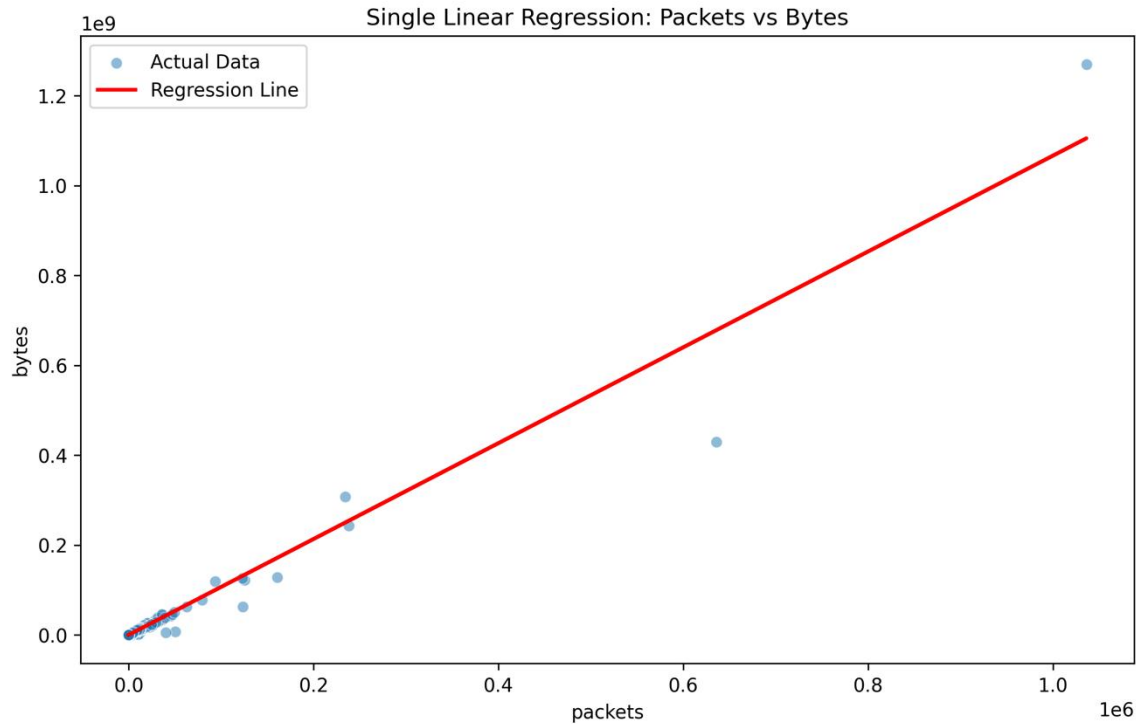Key Predictors: packets, elapsed_time_sec, nat_source_port, action

Firewall Actions: Allowed vs Blocked

Top 10 Target Ports (Traffic Type)

# Single Predictor Model



Single Linear Regression: Packets vs Bytes

Model: Bytes = β0 + β1·Packets

Packets chosen due to strong intuitive relationship

# Single Regression Results

Intercept ≈ -12,585

Slope ≈ 1,066

Each packet adds ~1066 bytes

p-value ≈ 0.000

T statistic ≈ 1109.02

$R^2$ ≈ 0.949

Residual Sum of Squares (RSS) ≈ $1.046 \times 10^{17}$

Residual Standard Error (RSE) ≈ 1 263 656 bytes ≈ 1.2 Megabytes
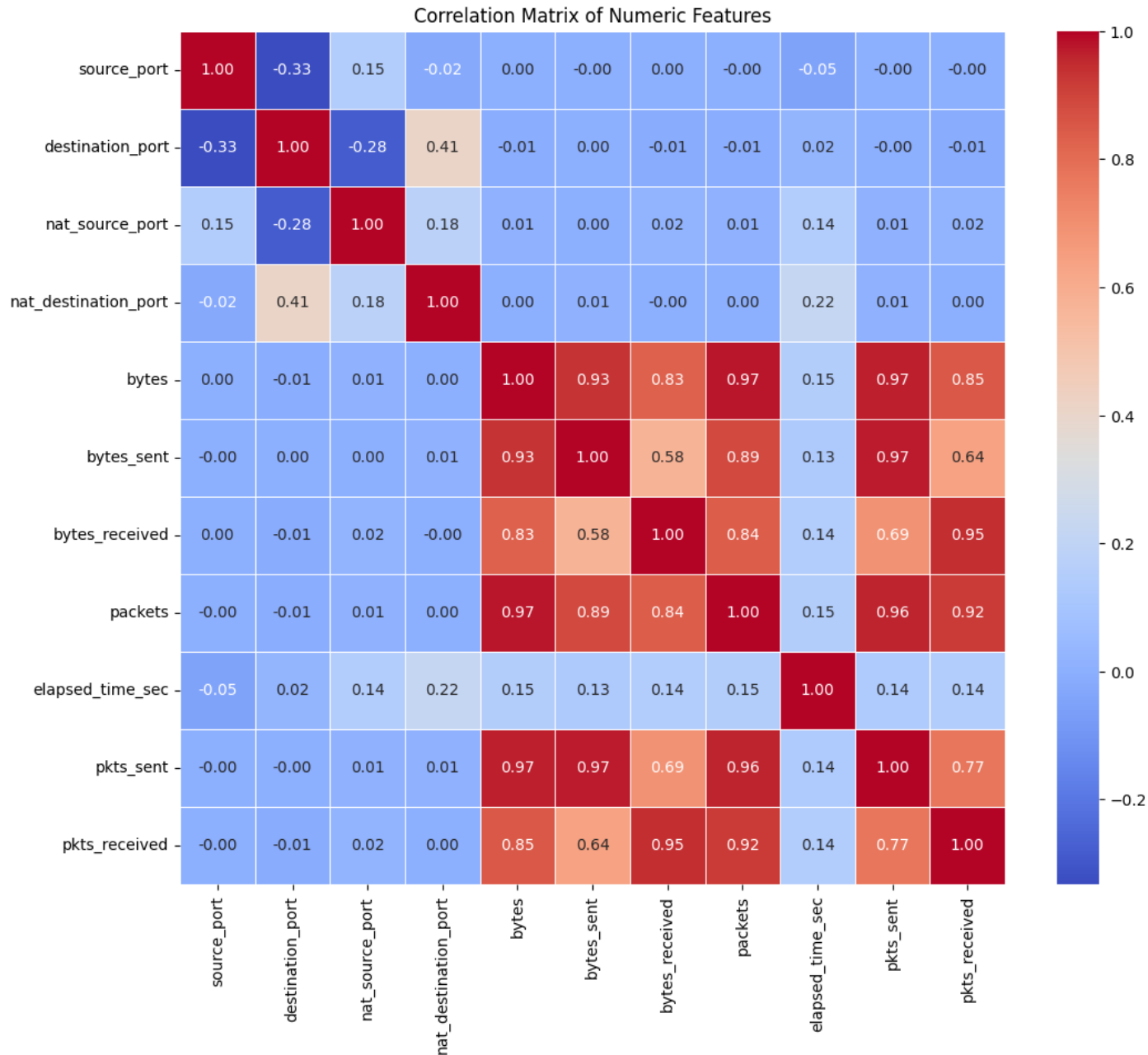
Correlation r ≈ 0.974

# Multiple Regression Model

Predictors added:

- Elapsed time (sec)

- NAT source port

| Metric | Single predictor | Multiple predictors | Conclusion |
|---|---|---|---|
| R-squared | 0.9494 | 0.9495 | Negligible improvement |
| P-values | < 0.05 | All < 0.05 | All variables are useful |

# Correlation Matrix



Correlation Matrix of Numeric Features

# Forward Selection

```
--- STARTING FORWARD SELECTION ---
Added: packets                | New Adj. R-squared: 0.949414
Added: bytes_sent             | New Adj. R-squared: 0.971599
Added: bytes_received         | New Adj. R-squared: 1.000000
(Stopping: Adding 'source_port' did not improve the model.)
-----------------------------------------------------
FINAL SUBSET SELECTED: ['packets', 'bytes_sent', 'bytes_received']
FINAL ADJ. R-SQUARED:  1.000000
```

# Prediction & Interaction Analysis

Packets=50, Time=10s, NAT Port=5000

Predicted Bytes ≈ 43,656

95% confidence interval [32,084, 55,228]

Packets × Action interaction

We tested if 'Action' (Allow vs Deny) affects the slope.

Result: Significant Interaction ($p < 0.05$).

'Allowed' packets carry more data (payload) than 'Denied' packets (headers only).

# Polynomial Regression

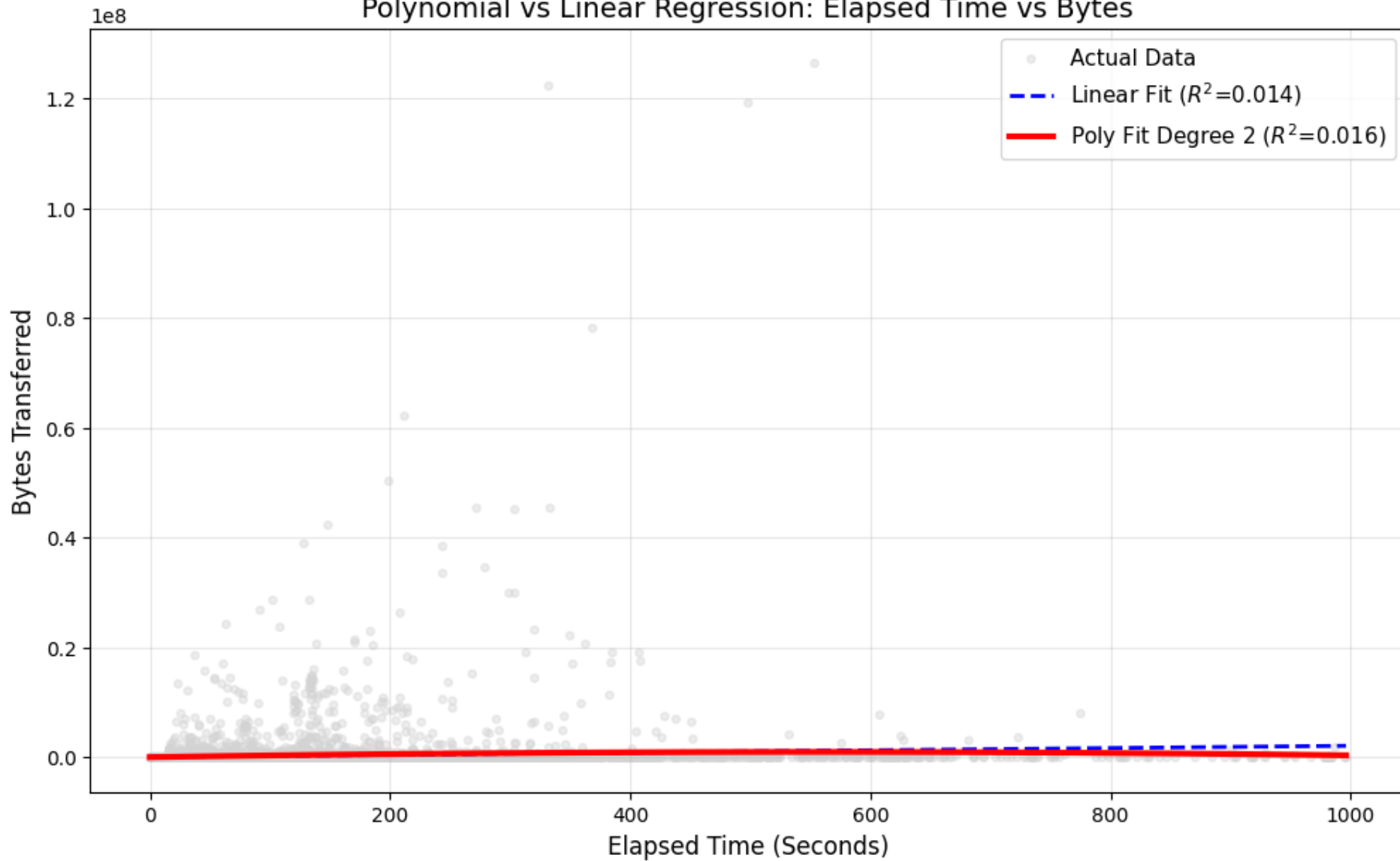Hypothesis: Does traffic grow non-linearly with time?

Model Comparison (Time vs Bytes):

1. Linear Fit: $R^2$ = 0.0136

2. Polynomial (Deg 2): $R^2$ = 0.0161

Polynomial fits slightly better, but Time is a weak predictor compared to Packets.

Polynomial vs Linear Regression: Elapsed Time vs Bytes

# Conclusion

Best Model: Single Predictor (Packets) is sufficient ($R^2 \approx 0.95$).

Adding Time and Port increased complexity but added negligible accuracy (< 0.001 improvement).

The volume of data transferred is almost entirely determined by the number of packets sent, regardless of connection duration.

Thank you for your attention