# Where Clause Column Selection using Decision Tree

Zihua Liu
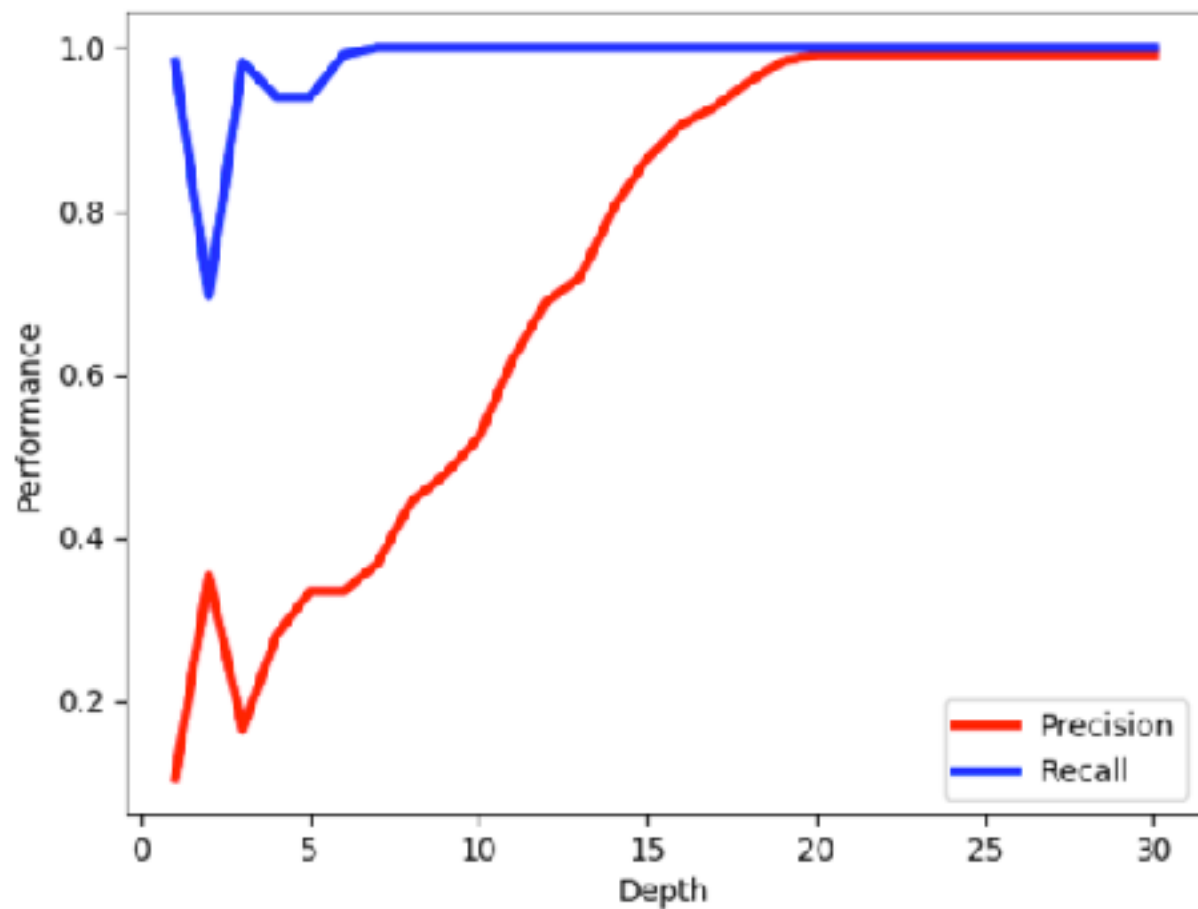
# Dataset

- Training Set

  - Total 2046 columns

  - 115 Positive

  - 1931 Negative

- Test Set

  - Total 559 Columns
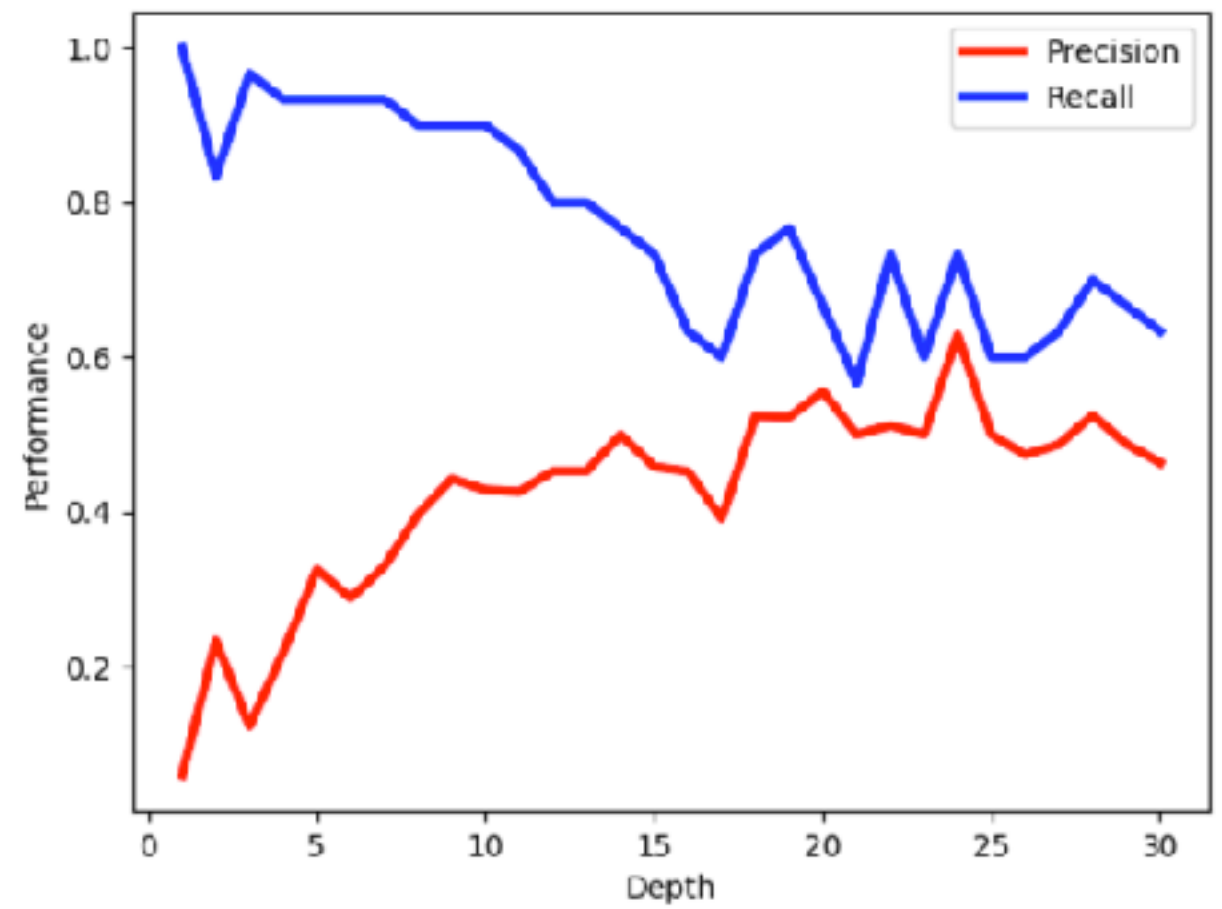
  - 30 Positive

  - 529 Negative

# Features

- Number of Columns

- Word2Vec proximities

- Column Types Probabilities

- Question Types (One hot)

- Min Edit Distance between question and column content

- Average Cell Length

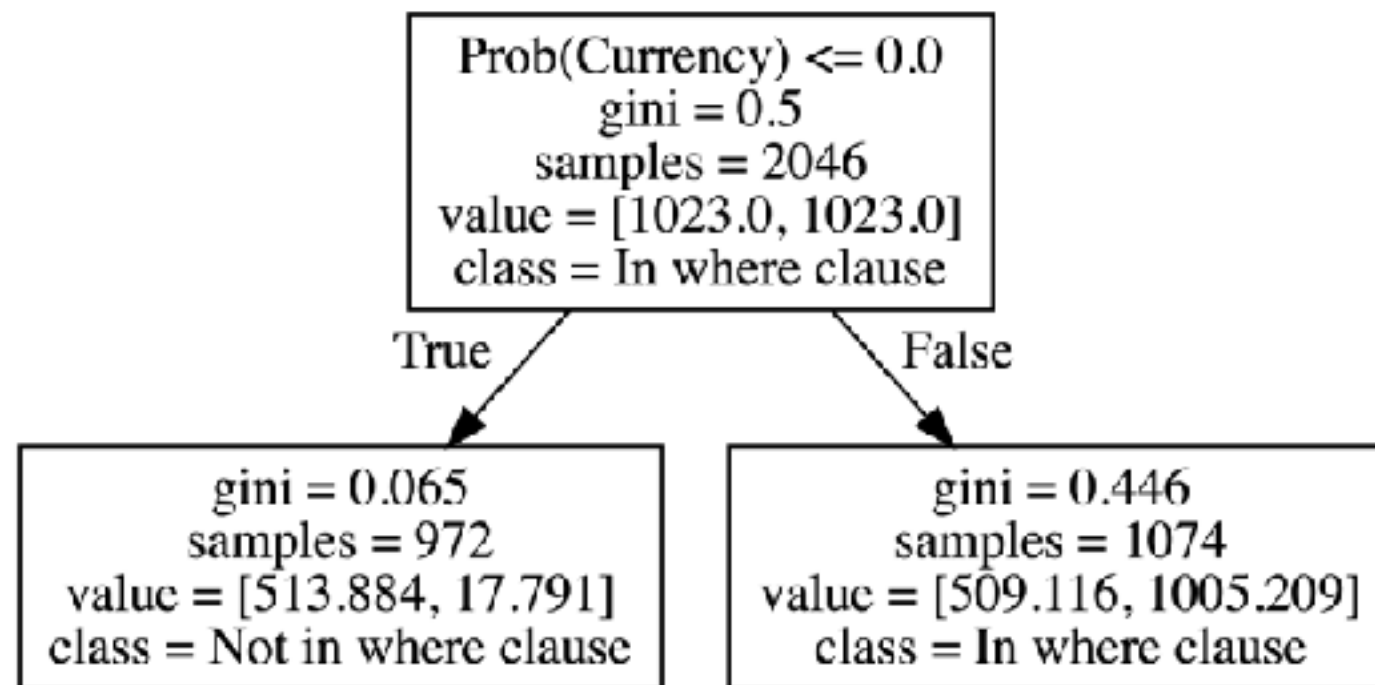- Number of Rows

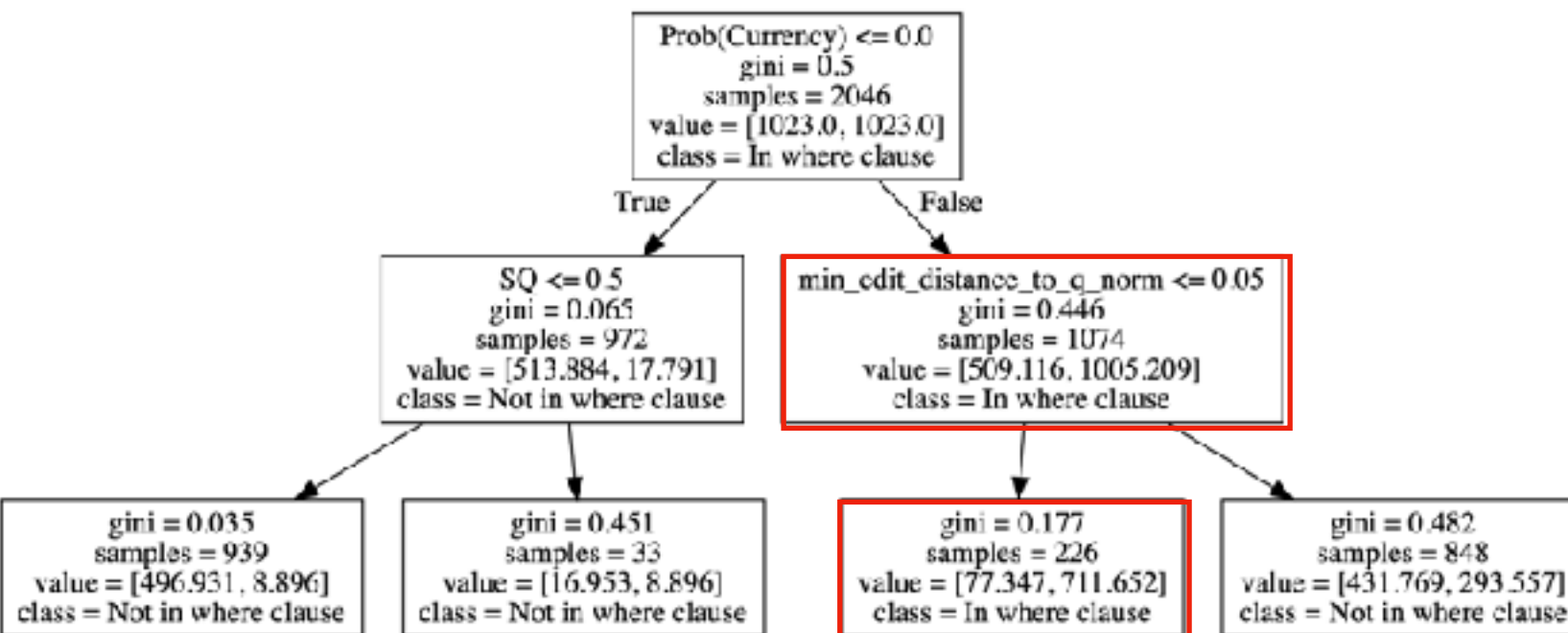- Whether in Select Columns

# Performance



**Training Set**

**Test Set**

# Depth = 1



| | Label = 1 | Label = 0 |
|---|---|---|
| Pred = 1 | 113 | 961 |
| Pred = 0 | 2 | 970 |

| | Label = 1 | Label = 0 |
|---|---|---|
| Pred = 1 | 30 | 474 |
| Pred = 0 | 0 | 55 |

# Depth = 2



| | Label = 1 | Label = 0 |
|---|---|---|
| **Pred = 1** | 80 | 146 |
| **Pred = 0** | 35 | 1785 |

| | Label = 1 | Label = 0 |
|---|---|---|
| **Pred = 1** | 25 | 82 |
| **Pred = 0** | 5 | 447 |

# Depth = 4

# Depth = 24

| | Label = 1 | Label = 0 |
|---|---|---|
| Pred = 1 | 115 | 1 |
| Pred = 0 | 0 | 1930 |

| | Label = 1 | Label = 0 |
|---|---|---|
| Pred = 1 | 22 | 13 |
| Pred = 0 | 8 | 516 |

# Thanks