# Column Projection

Zihua

# Feature Engineering: Column Types

- Text

- Numerical

- Currency

- Percentage

- Datetime

- Boolean

- URL

# Feature Engineering: Column Types

- Cons of Rule-Based Approach

  - Low Robustness

  - Overfitting on the training set

  - Feature vector is one-hot, while columns may contain multiple types of information
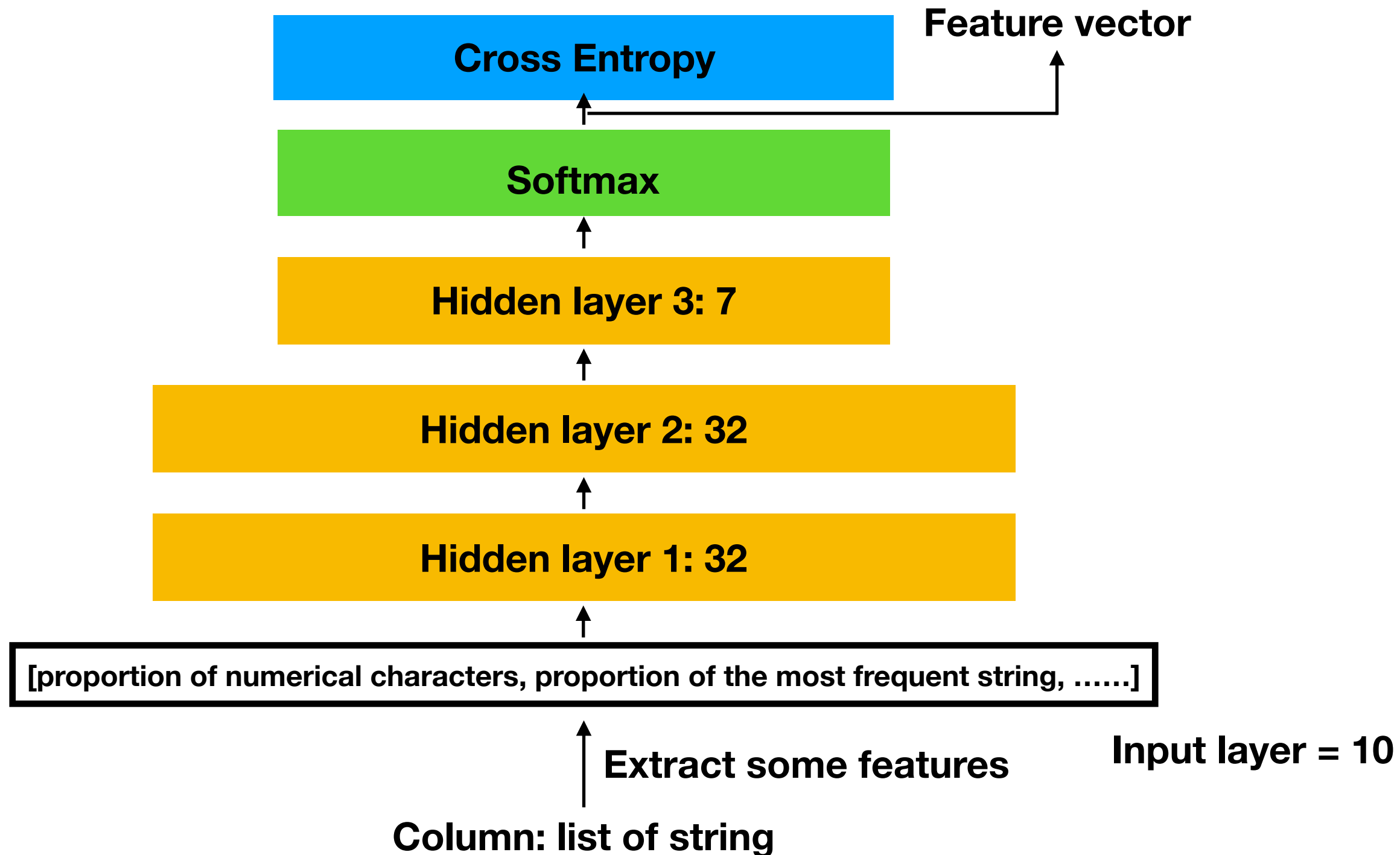
# Feature Engineering: Column Types

If use rule-based algorithm, then "Name" column should be attached with type "Text", so the datetime information in the column would be lost in the one hot vector.

| url | Number | Number_link | Presidency | Name | Prior_office | Party | Term | Vice_President |
|---|---|---|---|---|---|---|---|---|
| https://en.wikipedia.org/wiki/List_of_Presidents_of_the_United_States#Presidents | 1 | https://en.wikipedia.org/wiki/Presidency_of_George_Washington | April 30, 1789 [e] – March 4, 1797 | George Washington 1732–1799 (Lived: 67 years) [3][4][5] | Commander-in-Chief of the Continental Army (1775–1783) | Unaffiliated [2] | (1788–89) 1 (1789) | John Adams [f][g] |
| https://en.wikipedia.org/wiki/List_of_Presidents_of_the_United_States#Presidents | 1 | https://en.wikipedia.org/wiki/Presidency_of_George_Washington | April 30, 1789 [e] – March 4, 1797 | George Washington 1732–1799 (Lived: 67 years) [3][4][5] | Commander-in-Chief of the Continental Army (1775–1783) | | (1792) 2 (1793) | John Adams [f][g] |
| https://en.wikipedia.org/wiki/List_of_Presidents_of_the_United_States#Presidents | 2 | https://en.wikipedia.org/wiki/Presidency_of_John_Adams | March 4, 1797 – March 4, 1801 | John Adams 1735–1826 (Lived: 90 years) [6][7][8] | 1st Vice President of the United States; 1st Vice President of the United States | Federalist | (1796) 3 (1797) | Thomas Jefferson [h] |
| https://en.wikipedia.org/wiki/List_of_Presidents_of_the_United_States#Presidents | 3 | https://en.wikipedia.org/wiki/Presidency_of_Thomas_Jefferson | March 4, 1801 – March 4, 1809 | Thomas Jefferson 1743–1826 (Lived: 83 years) [9][10][11] | 2nd Vice President of the United States; 2nd Vice President of the United States | Democratic-Republican | (1800) 4 (1801) | Aaron Burr March 4, 1801 – March 4, 1805 |

# Feature Engineering: Column Types

- Machine Learning approach

  - Manually labeled all the columns in the training set

  - Trained a MLP for classification

  - Generate a probability distribution over all types for each column

# Feature Engineering: Column Types

# Feature Engineering: Column Types



**Accuracy: 93 % on training set**

| Column Name | P(DateTime) | P(Text) | P(Currency) | P(Boolean) | P(Numerical) | P(Percentage) | P(URL) |
| --- | --- | --- | --- | --- | --- | --- | --- |
| Name | 0.58551943 | 0.3745269 | 0.000408226 | 0.004346717 | 0.03398786 | 0.001172554 | 3.83E-05 |
| Party | 0.001754907 | 0.9906969 | 3.58E-05 | 0.000395942 | 0.007047001 | 2.46E-05 | 4.19E-05 |
| Term | 0.8488813 | 0.034917403 | 0.000619899 | 0.007573995 | 0.10370393 | 0.004300023 | 3.36E-06 |

# Column projection

- Number of columns (Lucile)

- Word2vec proximities (Lucile)

- Column Types (Zihua)

- Question Types (Nithin)—> Transform to one-hot vector

# Column projection

For training set, there are total 2046 cases, 273 positive (label = 1), 1773 negative (label = 0), so copy positive cases 6 times

SoftmaxCrossEntropy

Hidden layer 4: 2

Batch Size = 64

BatchNorm layer

Hidden layer 3: 8

BatchNorm layer

Hidden layer 2: 16

BatchNorm layer

Hidden layer 1: 32

concatenated feature vector

size = 19

# Column projection

## Training set

|       | Postitive | Negative |
|-------|-----------|----------|
| True  | 222       | 1407     |
| False | 51        | 366      |
|       | 273       | 1773     |

## Test set

|       | Postitive | Negative |
|-------|-----------|----------|
| True  | 49        | 394      |
| False | 17        | 99       |
|       | 66        | 493      |