# WEBSITE CATEGORIZATION

CAPSTONE PROJECT PRESENTATION

# OUTLINE

# I. ABOUT OUR PROJECT

CATEGORIZATION CAN BE DEFINED AS "THE PROCESS BY WHICH INDIVIDUALS GROUP OBJECTS OR EVENTS INTO CATEGORIES".

**WEBSITE CATEGORIZATION:**

- **INPUT** INFORMATION IS THE CONTENT OF WEBSITES IN TEXT FORM.
- **OUTPUT** LABEL FOR EACH CONTENT

- Source list of websites: https://www.similarweb.com
- **11 categories**
- Average **50–60 websites** per topic

```
1  const categoriesSimilarWeb = [
2    "adult",
3    "arts-and-entertainment",
4    "business-and-consumer-services",
5    "computers-electronics-and-technology",
6    'sports',
7    'science-and-education',
8    'food-and-drink',
9    'travel-and-tourism',
10   'health',
11   'pets-and-animals',
12   'law-and-government'
13  ]
```

We only select websites from nations where English is widely spoken:

```
const region = ["united-kingdom", "united-states", "australia",
                "canada", "ireland", "new-zealand", "netherlands"];
```

similarweb ⊢ --- → AXIOS --- → Categorized Website List

Save to Json file ⊢ --- ⊣ Manual Filtering

**data.csv**

AXIOS

English words ← --- CHEERIO --- ⊣ HTML text

# II. DATA SCIENCE METHODOLOGY
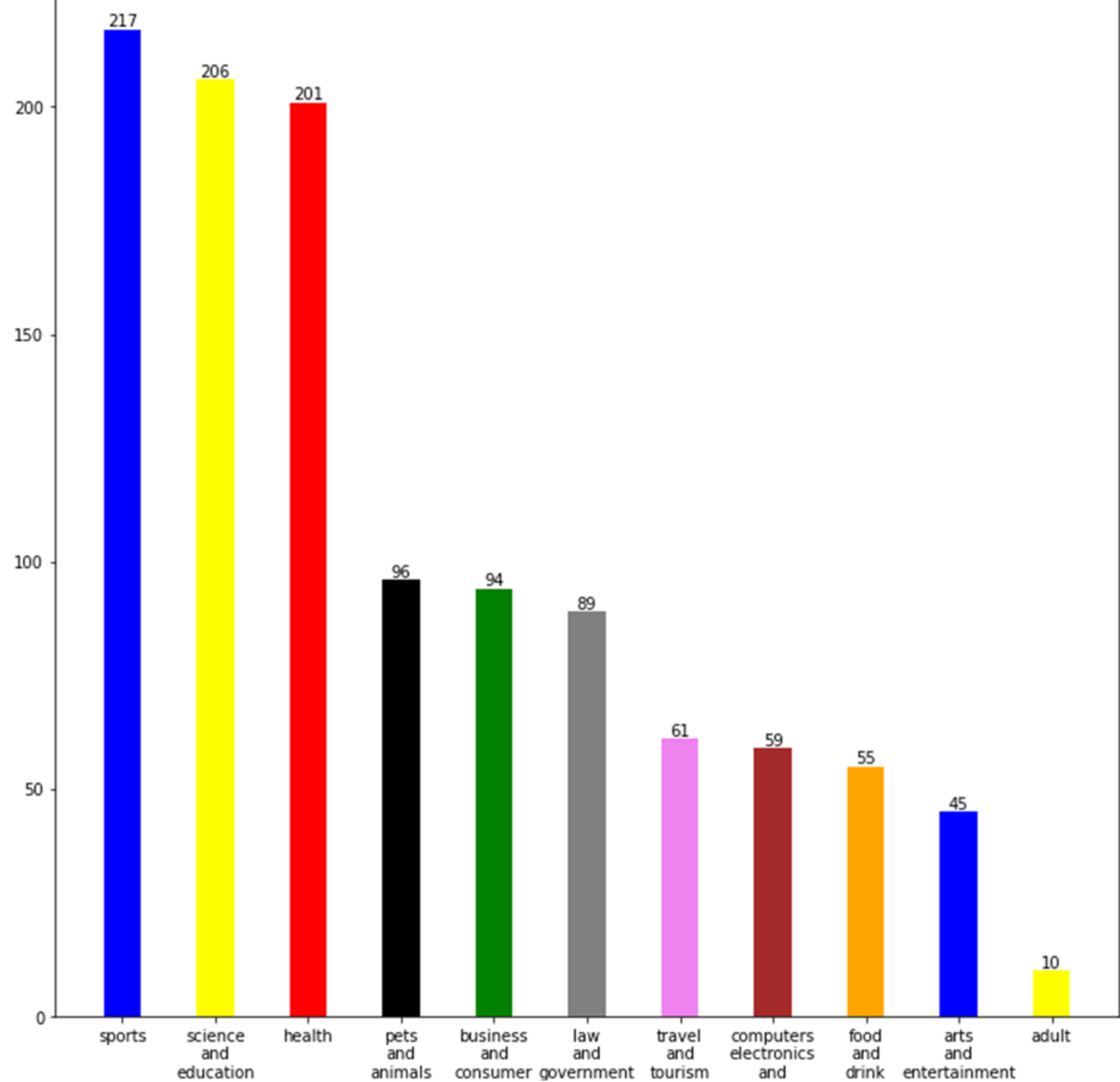## 1. DATA CRAWLING

## Note:

- Only use **Javascript library** to crawl website for convenience

- Only crawl **landing page** or **homepage** of each website

- We **don't handle Javascript** implementation and only focus on **SSR website**, which has plenty of HTML text (> 1000 English words)

- Remove noise: Manual filtering to **remove unreadable and non-English** content.

*LOWERING ALL THE TEXT

*REMOVE ALL HYPERLINK , URL IN ALL DOCUMENTS

*REMOVE ALL PUNCTUATION SUCH AS ";" , "?" ,"!" AND OTHER UNNECESSARY CHARACTER OR SYMBOL

*REMOVE ALL HTML TAG

*REMOVE ALL STOP WORD SUCH AS A , AN ,THE ,….

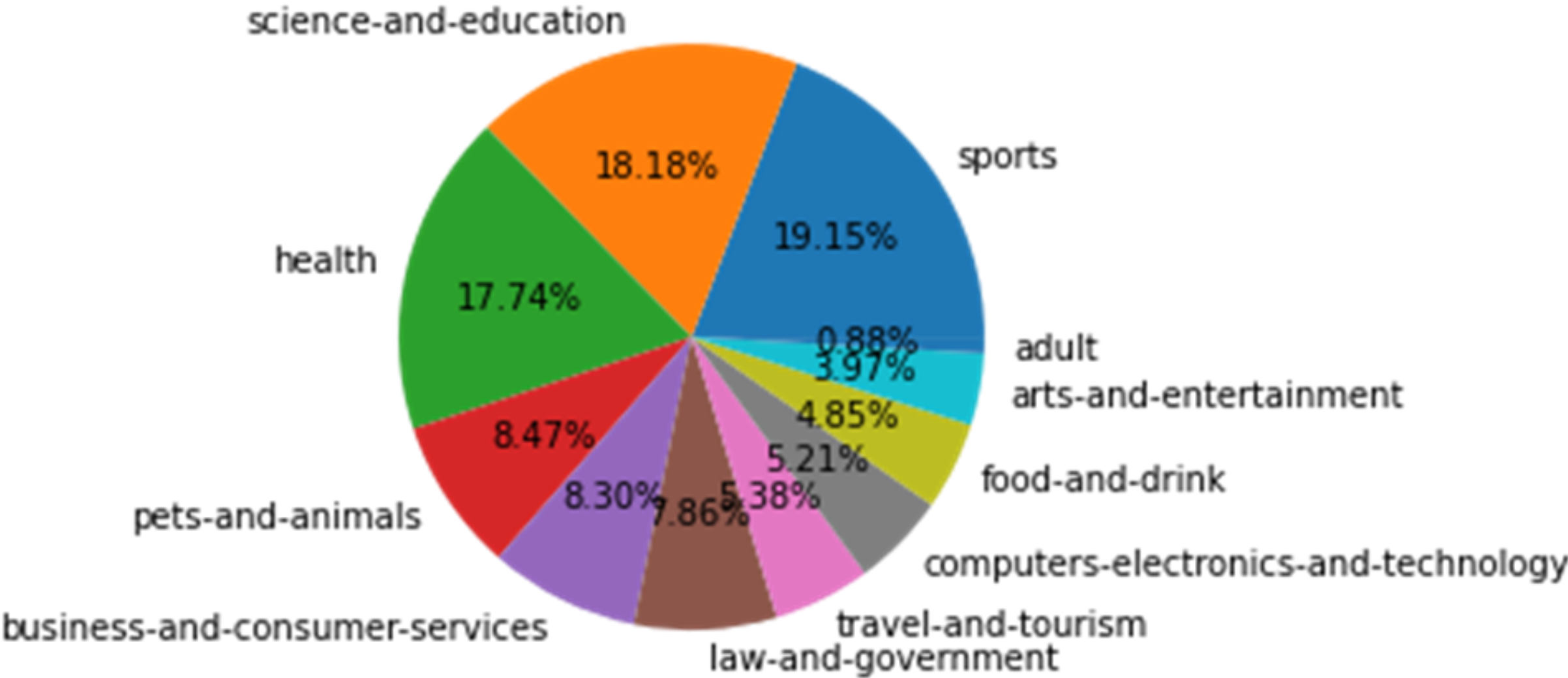*USING SNOWBALLSTEMMER TO STEM THE ALL THE WORD

# 3. DATA VISUALIZATION

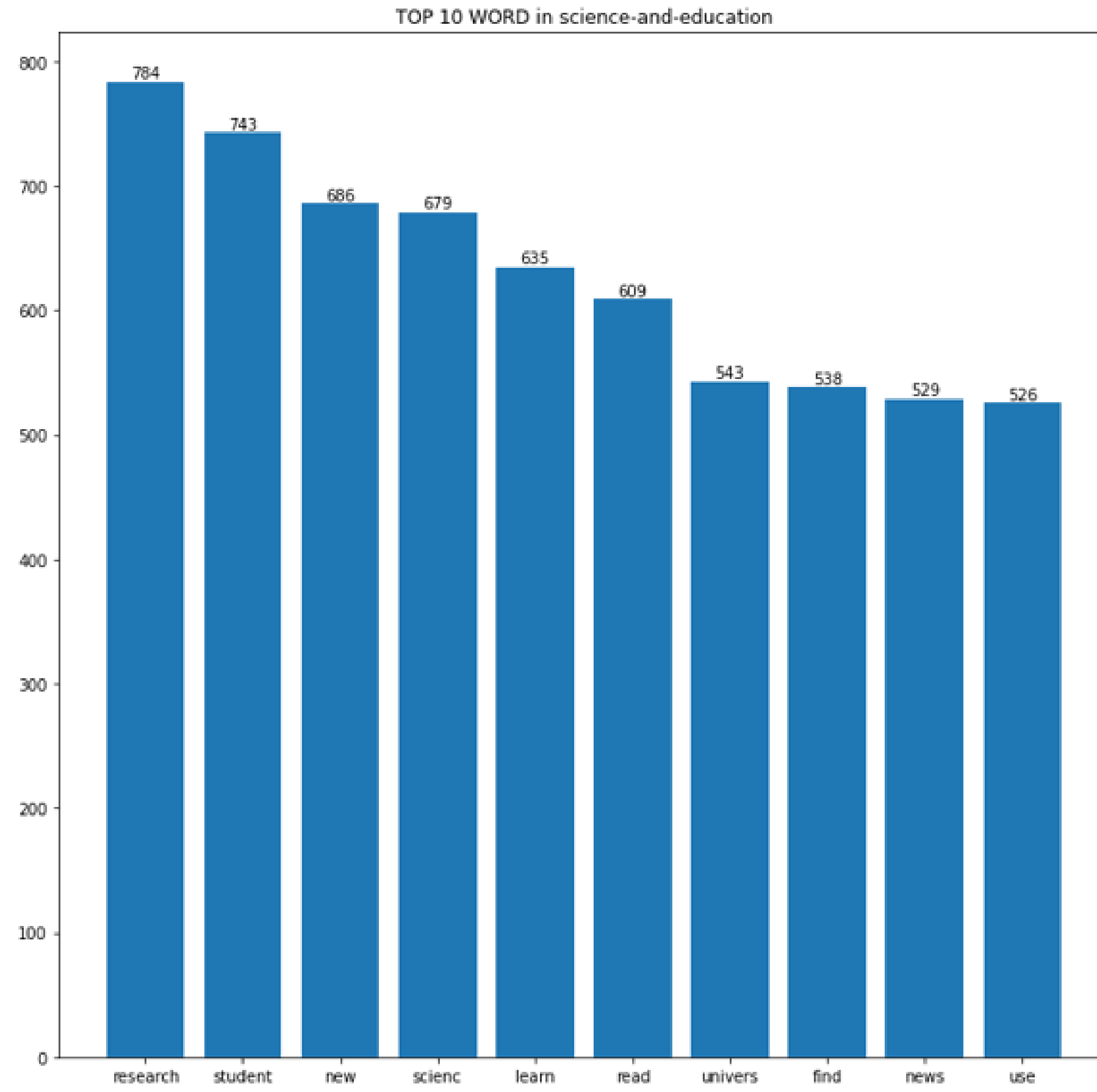The bar chart illustrates the amount of crawled websites.

# 3. DATA VISUALIZATION

The pie chart representing the percentage of the
number of websites in each topic

# 3. DATA VISUALIZATION

# 3. DATA VISUALIZATION



TOP 10 WORD in science-and-education

# 4. FEATURE EXTRACTION

**TF(TERM FREQUENCY):**

$$TF = \frac{\text{Number of times word appear in a document}}{\text{Total number of word in that document}}$$
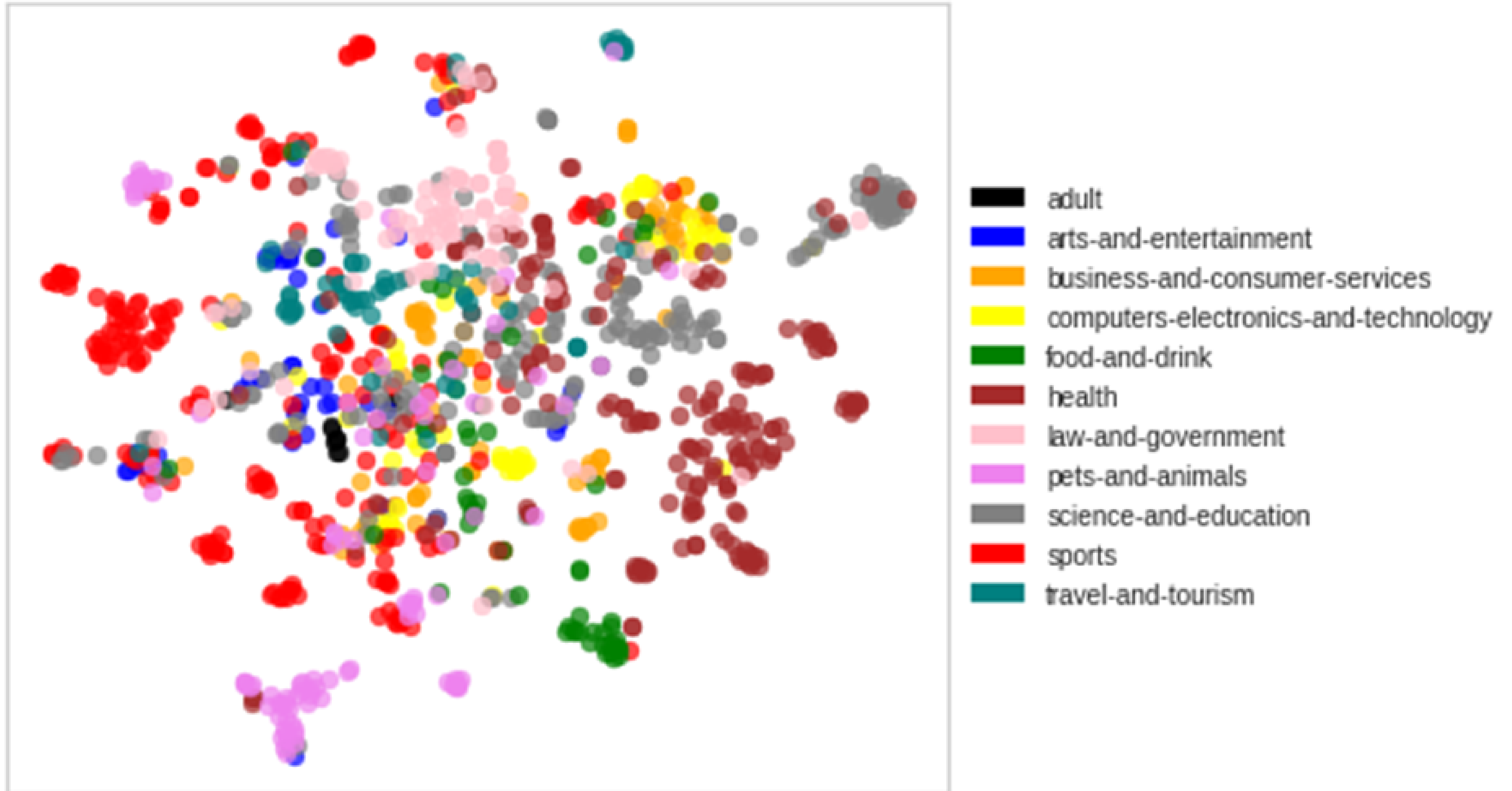
## IDF(INVERSE DOCUMENT FREQUENCY )

$$IDF = \log\left(\frac{\text{Total number of documents}}{\text{Document that contain word W}}\right)$$

### => TF-IDF SCORE= TF*IDF

## TF-IDF VECTORIZER:

```python
tfidf = TfidfVectorizer(min_df=0.01,max_df=0.85,max_features = 1500,ngram_range=(1,3))
docs = tfidf.fit_transform(df.clean_text)
```
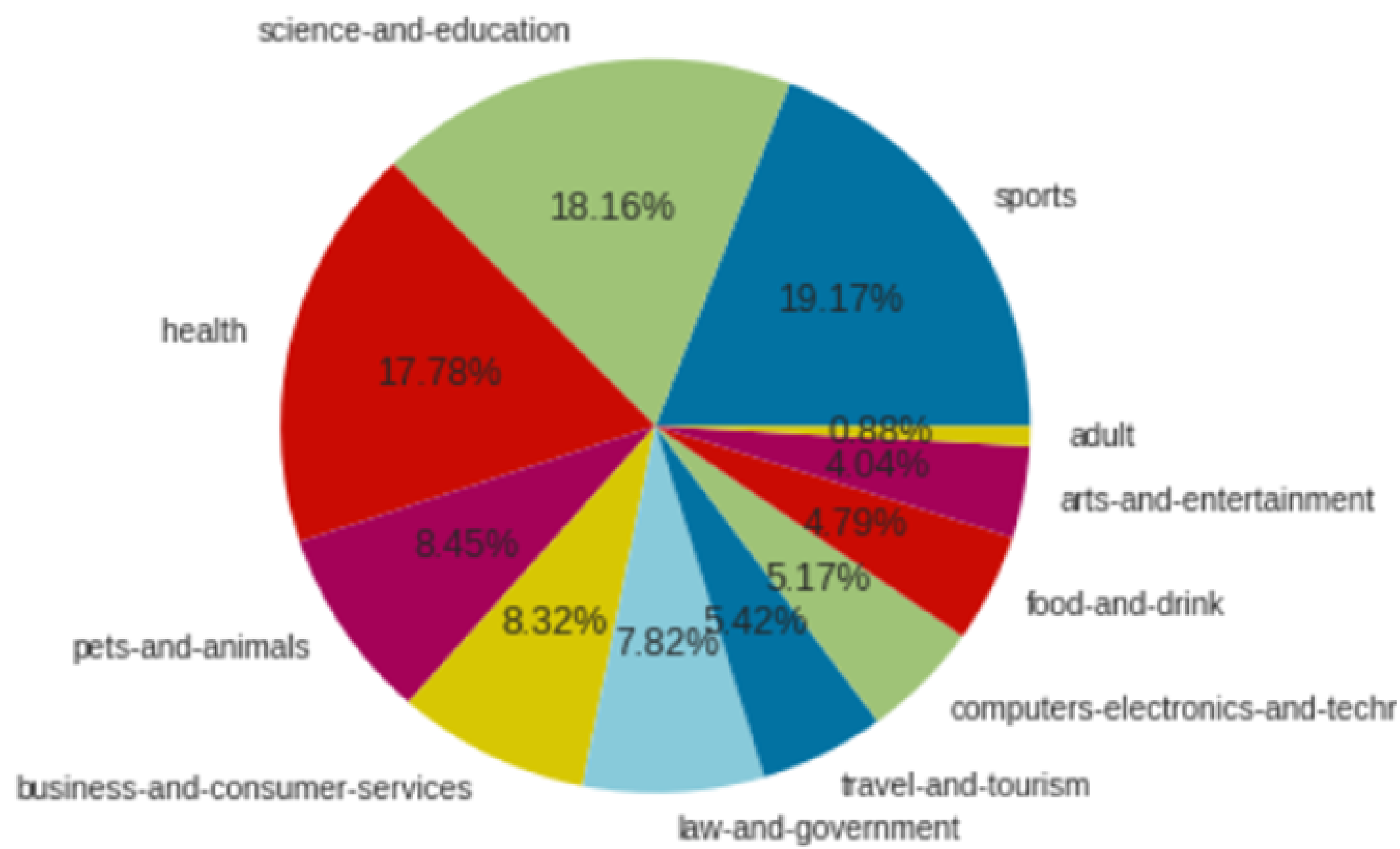
TSNE Projection of 1133 Documents

| Color | Label |
|-------|-------|
| black | adult |
| blue | arts-and-entertainment |
| orange | business-and-consumer-services |
| yellow | computers-electronics-and-technology |
| green | food-and-drink |
| dark red | health |
| pink | law-and-government |
| violet | pets-and-animals |
| gray | science-and-education |
| red | sports |
| teal | travel-and-tourism |

# 4. MODEL TRAINING

Divide the data set into 2 subsets:

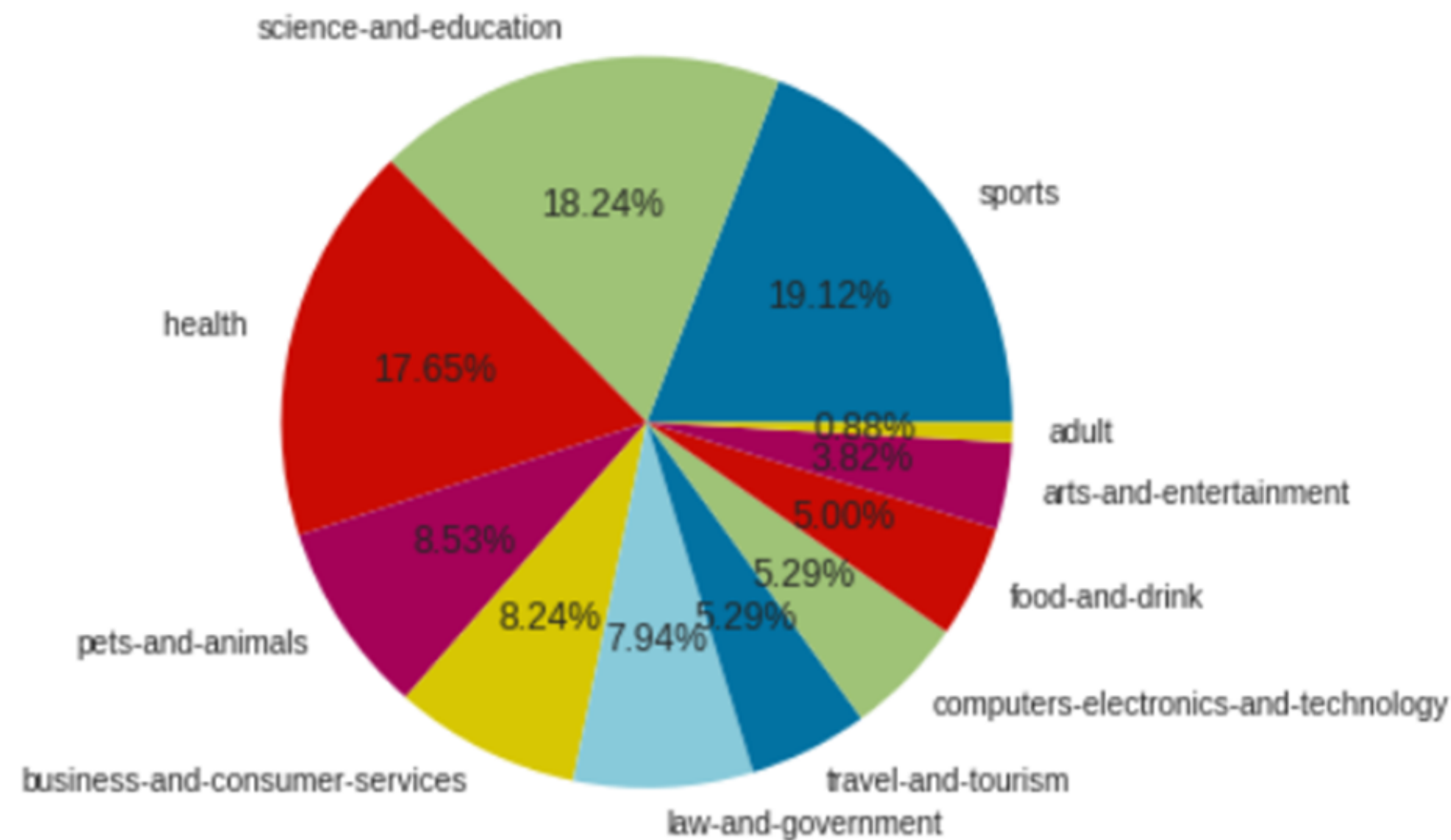- Train set: 70%

- Test set: 30 %

# 4. MODEL TRAINING

The percentage of websites crawled by topic in training set and test set:



**Training set(793 rows)**

# 4. MODEL TRAINING

The percentage of websites crawled by topic in training
set and test set:



**Test set(340 Rows)**
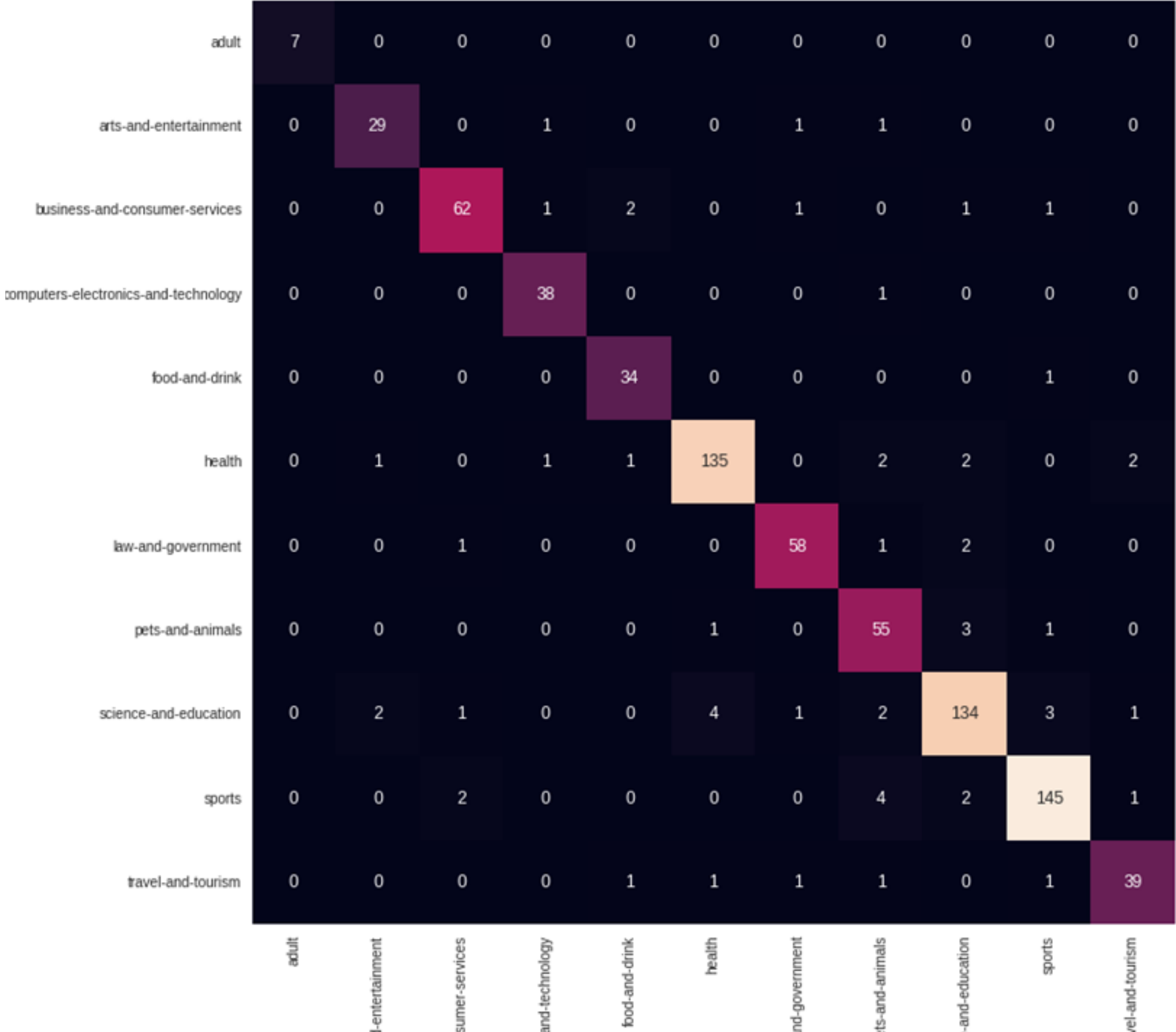
# 4. MODEL TRAINING

**NAIVE BAYES:**

$$P(A|B) = \frac{P(B|A)\, P(A)}{P(B)}$$

**THE ATTRIBUTES ARE CONDITIONALLY INDEPENDENT GIVEN CLASSIFICATION**

$$c_{MAP} = \arg\max_{c_i \in C} P(c_i) . \prod_{j=1}^{n} P(z_j | c_i)$$

# II. RESULT

## 1. FINAL RESULT



*TRAINING:

-ACCURACY:93%

# II. RESULT

## 1. FINAL RESULT



*TEST

ACCURACY:77%

# II. RESULT

## 2. DIFFICULTY

*CRAWL WEB

*MODEL ACCURACY

# II. RESULT
## 3. FUTURE IMPROVEMENT

*FIND AND CRAWL MORE WEB TO INCREASE  DATASET

*HANDLE MULTI-LABLE PROBLEM

*TRY ANOTHER MODEL

# THANKS FOR LISTENING OUR PRESENTATION