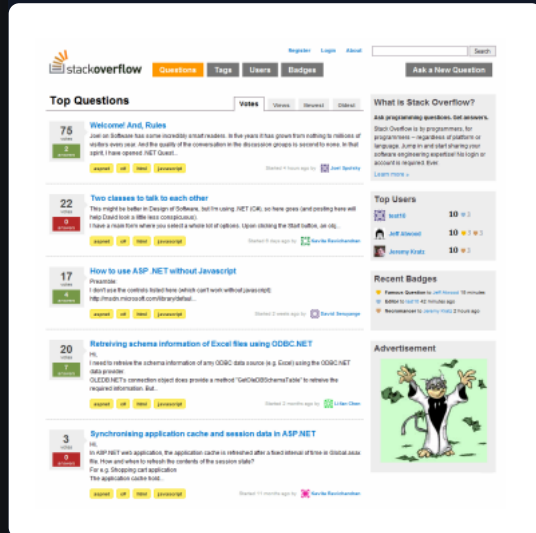# Stack Exchange Question Classifier ★

Problem | Submissions | Leaderboard | Editorial 🔒

Stack Exchange is an information powerhouse, built on the power of crowdsourcing. It has 105 different topics and each topic has a library of questions which have been asked and answered by knowledgeable members of the StackExchange community. The topics are as diverse as travel, cooking, programming, engineering and photography.



We have hand-picked ten different topics (such as Electronics, Mathematics, Photography etc.) from Stack Exchange, and we provide you with a set of questions from these topics.

Given a question and an excerpt, your task is to identify which among the 10 topics it belongs to.

**Getting started with text classification**

For those getting started with this fascinating domain of text classification, here's a wonderful Youtube video of Professor Dan Jurafsky from Stanford, explaining the Naive Bayes classification algorithm, which you could consider using as a starting point.

## Input Format

The first line will be an integer N. N lines follow each line being a valid JSON object. The following fields of raw data are given in json

- question (string) : The text in the title of the question.
- excerpt (string) : Excerpt of the question body.
- topic (string) : The topic under which the question was posted.

The input for the program has all the fields but topic which you have to predict as the answer.

## Constraints

1 <= N <= 22000

topic is of ascii format

question is of UTF-8 format

excerpt is of UTF-8 format

## Output Format

For each question that is given as a JSON object, output the topic of the question as predicted by your model separated by newlines.

The training file is available here. It is also present in the current directory in which your code is executed.

## Sample Input

```
12345
json_object
json_object
json_object
.
.
.
json_object
```

## Sample Output

```
electronics
security
photo
.
.
.
mathematica
```

Sample testcases can be downloaded here for offline training. When you submit your solution to us, you can assume that the training file can be accessed by reading "training.json" which will be placed in the same folder as the one in which your program is being executed.

## Scoring

While the contest is going on, the score shown to you will be on the basis of the Sample Test file. The final score will be based on the Hidden Testcase only and there will be no weightage for your score on the Sample Test.

Score = MaxScore for the test case * (C/T)

Where C = Number of topics identified correctly and

T = total number of test JSONs in the input file.

Change Theme    Language    Python 3

```python
5    import sys
6    from sklearn.feature_extraction.text import CountVectorizer
7    from sklearn.naive_bayes import MultinomialNB
8    from sklearn.svm import LinearSVC
9    from sklearn.model_selection import train_test_split
10
11   file_name = 'training.json'
12   with codecs.open(file_name, 'r', encoding='iso-8859-1') as f:
13       data = f.read()
14   data = data.split("\n")
15
16   N = int(data[0].strip())
17   topics = []
```

```
18    questions = ·[]
19    excerpts = ·[]
20    rows = ·[]
21    for·i·in·range(1,N+1):
22    ····item = ·json.loads(data[i])
23    ····topics.append(item['topic'])
24    ····questions.append(item['question'])
25    ····excerpts.append(item['excerpt'])
26
27    ····row = ·[item['topic'], ·item['question'], ·item['excerpt']]
28    ····rows.append(row)
29    colnames = ·["topic", ·"question", ·"excerpt"]
30    training_df = ·pd.DataFrame(rows, ·columns= ·colnames)
31
32
33    df = ·sys.stdin.readlines()
34    df = ·[line.rstrip() ·for·line·in·df]
35    N = ·int(df[0].strip())
36    topics = ·[]
37    questions = ·[]
38    excerpts = ·[]
39    rows = ·[]
40    for·i·in·range(1,N+1):
41    ····item = ·json.loads(df[i])
42    ····
43    ····questions.append(item['question'])
44    ····excerpts.append(item['excerpt'])
```

Line: 67 Col: 13

⬆ Upload Code as File          ☐ Test against custom input          Run Code          Submit Code

# Congratulations!

You have passed the sample test cases. Click the submit button to run your code against all the test cases.

✅ **Sample Test case 0**

Compiler Message

**Success**

Input (stdin)                                                                    Download

```
1    21345

2    {"question":"Are there any SciFi treatments of time travel that avoid the typical
     paradoxes? [duplicate]","excerpt":"Possible Duplicate:\n  Why do time-travel
     stories often have the characters "returning" to the future?  \n\n\n\n\nThe
     possibility of time travel normally creates paradoxes. If you can travel into the
     ...\r\n        "}

3    {"question":"How to auto strip hyperlinks &amp; images in wordpress
     post","excerpt":"I creat a post form on front-end for wordpress members by use
     DJD Site Post plugin.\nPlugin url: http://wordpress.org/extend/plugins/djd-site-
     post/\n\nAnd now i want to auto strip all hyperlink &amp; ...\r\n        "}

4    {"question":"Why do fantasy writers depict pointy hats as the headgear of choice
```