
Build a Question Answering Engine in Minutes

@Milvus.io

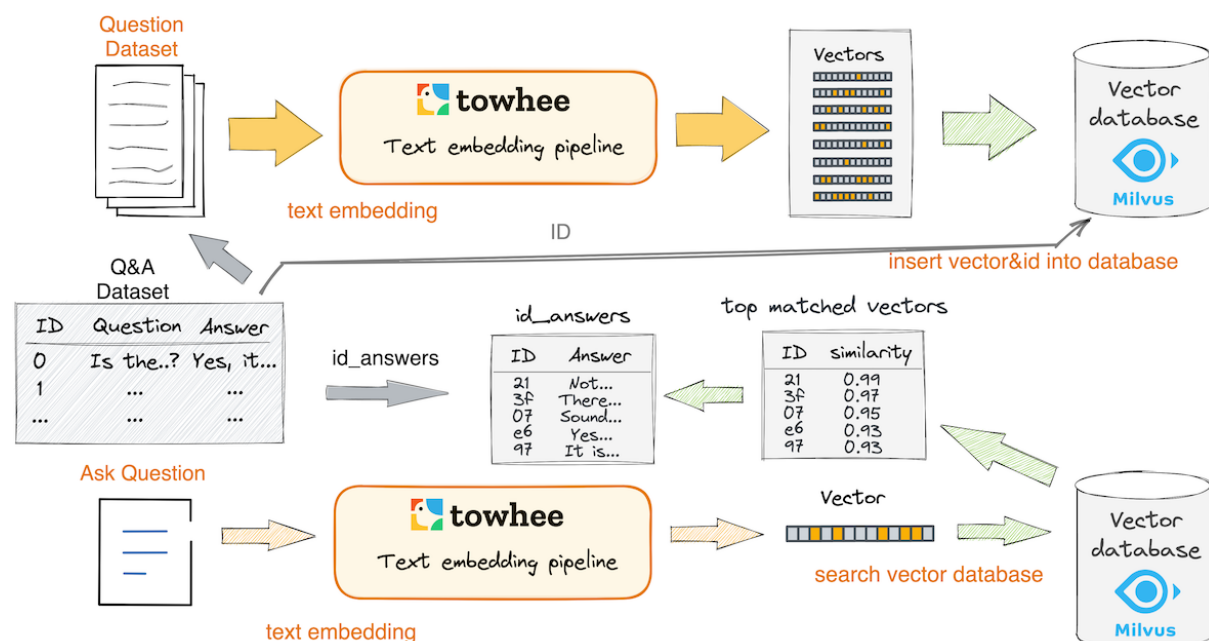
6/9/2022

Build a Question Answering Engine in Minutes

Introduction

Question answering is a classic problem in the field of natural language processing. While it sounds like an easy problem to solve, there is still a lot of research going on to improve the techniques that we have now. A large part of solving questions is finding questions that are similar to the one being asked.

This example will show you how to find the similar asked question and get the answer. The basic idea behind question answering is to use Towhee to generate embedding from the question dataset and compare the input question with the embedding stored in Milvus.



Preparations

Install Dependencies

First we need to install dependencies such as pymilvus, towhee and gradio.

```
1 $ python -m pip install -q pymilvus towhee gradio
```

Prepare the Data

There is a subset of the InsuranceQA Corpus (1000 pairs of questions and answers) used in this demo, everyone can download on Github.

```
1 $ curl -L https://github.com/towhee-io/examples/releases/download/data/question_answer.csv -O
```

question_answer.csv: a file containing question and the answer.

Let's take a quick look:

```
1 import pandas as pd
2
3 df = pd.read_csv('question_answer.csv')
4 df.head()
```

	id	question	answer
0	0	Is Disability Insurance Required By Law?	Not generally. There are five states that requ...
1	1	Can Creditors Take Life Insurance After ...	If the person who passed away was the one with...
2	2	Does Travelers Insurance Have Renters Ins...	One of the insurance carriers I represent is T...
3	3	Can I Drive A New Car Home Without Ins...	Most auto dealers will not let you drive the c...
4	4	Is The Cash Surrender Value Of Life Ins...	Cash surrender value comes only with Whole Lif...

To use the dataset to get answers, let's first define the dictionary:

- **id_answer**: a dictionary of id and corresponding answer

```
1 id_answer = df.set_index('id')['answer'].to_dict()
```

Create Milvus Collection

Before getting started, please make sure you have installed milvus. Next to define the function `create_milvus_collection` to create collection in Milvus that uses the L2 distance metric and an IVF_FLAT index.

```
1 from pymilvus import connections, FieldSchema, CollectionSchema,
   DataType, Collection, utility
2
3 connections.connect(host='127.0.0.1', port='19530')
4
5 def create_milvus_collection(collection_name, dim):
6     if utility.has_collection(collection_name):
7         utility.drop_collection(collection_name)
```

```
8
9     fields = [
10     FieldSchema(name='id', dtype=DataType.INT64, description='ids',
11                 is_primary=True, auto_id=False),
12     FieldSchema(name='embedding', dtype=DataType.FLOAT_VECTOR,
13                 description='embedding vectors', dim=dim)
14     ]
15     schema = CollectionSchema(fields=fields, description='reverse image
16         search')
17     collection = Collection(name=collection_name, schema=schema)
18
19     # create IVF_FLAT index for collection.
20     index_params = {
21         'metric_type': 'L2',
22         'index_type': "IVF_FLAT",
23         'params': {"nlist": 2048}
24     }
25     collection.create_index(field_name="embedding", index_params=
26         index_params)
27     return collection
```

Load question embedding into Milvus

We first generate embedding from question text with dpr operator and insert the embedding into Milvus. Towhee provides a method-chaining style API so that users can assemble a data processing pipeline with operators.

```
1 import towhee
2
3 collection = create_milvus_collection('question_answer', 768)
4
5 dc = (
6     towhee.read_csv('question_answer.csv')
7     .runas_op['id', 'id'](func=lambda x: int(x))
8     .text_embedding.dpr['question', 'vec'](model_name="facebook/dpr-
9         ctx_encoder-single-nq-base")
10    .runas_op['vec', 'vec'](func=lambda x: x.squeeze(0))
11    .tensor_normalize['vec', 'vec']()
12    .to_milvus['id', 'vec'](collection=collection, batch=100)
13 )
14 print('Total number of inserted data is {}'.format(collection.
15     num_entities))
```

Explanation of Data Processing Pipeline Here is detailed explanation for each line of the code:

```
towhee.read_csv('question_answer.csv'): read tabular data from the file (id, question and answer columns);

.runas_op['id', 'id'](func=lambda x: int(x)): for each row from the data, convert the data type of the column id from str to int;

.text_embedding.dpr['question', 'vec'](model_name="facebook/dpr-ctx_encoder-single-nq-base"): use the facebook/dpr-ctx_encoder-single-nq-base model to generate the question embedding vector with the dpr operator in towhee hub;

.runas_op['vec', 'vec'](func=lambda x: x.squeeze(0)): the vec shape after dpr operator is (1, 768), so we need to squeeze it;

.tensor_normalize['vec', 'vec'](): normalize the embedding vector;

.to_milvus['id', 'vec'](collection=collection, batch=100): insert question embedding vector into Milvus;
```

Ask Question with Milvus and Towhee

Now that embedding for question dataset have been inserted into Milvus, we can ask question with Milvus and Towhee. Again, we use Towhee to load the input question, compute a embedding, and use it as a query in Milvus. Because Milvus only outputs IDs and distance values, we provide the `id_answers` dictionary to get the answers based on IDs and display.

```
1 dc = ( towhee.dc(['Is Disability Insurance Required By Law?'])
2       .text_embedding.dpr(model_name="facebook/dpr-ctx_encoder-single-
3         nq-base")
4       .runas_op(func=lambda x: x.squeeze(0))
5       .tensor_normalize()
6       .milvus_search(collection='question_answer', limit=1)
7       .runas_op(func=lambda res: [id_answer[x.id] for x in res])
8       .to_list()
```

Then we can get the answer about 'Is Disability Insurance Required By Law?'.

```
1 dc[0]
```

'Not generally. There are five states that require most all employers carry short term disability insurance on their employees. These states are: California, Hawaii, New Jersey, New York, and Rhode Island. Besides this mandatory short term disability law, there is no other legislative imperative for someone to purchase or be covered by disability insurance.'

Release a Showcase

We've done an excellent job on the core functionality of our question answering engine. Now it's time to build a showcase with interface. Gradio is a great tool for building demos. With Gradio, we simply need to wrap the data processing pipeline via a `chat` function:

```
1 def chat(message, history):
2     history = history or []
3     with towhee.api() as api:
4         qa_function = (
5             api.text_embedding.dpr(model_name="facebook/dpr-
6                 ctx_encoder-single-nq-base")
7                 .runas_op(func=lambda x: x.squeeze(0))
8                 .tensor_normalize()
9                 .milvus_search(collection='question_answer', limit
10                     =3)
11                 .runas_op(func=lambda res: [id_answer[x.id]+'\\n'
12                     for x in res])
13                 .as_function()
14         )
15     response = qa_function(message)[0]
16     history.append((message, response))
17     return history, history
18
19 import gradio
20
21 chatbot = gradio.Chatbot(color_map=("green", "gray"))
22 interface = gradio.Interface(
23     chat,
24     ["text", "state"],
25     [chatbot, "state"],
26     allow_screenshot=False,
27     allow_flagging="never",
28 )
29 interface.launch(inline=True, share=True)
```