

Cady Baltz
CS 4395.001 Human Language Technologies
Portfolio Assignment: Finding or Building a Corpus
3/11/2023

Introduction:

In this project, I designed a web crawler to find websites related to tourism in Japan. The goal of this program was to create a knowledge base of information that could be used to create a travel agent chatbot that can help answer questions about visiting Japan.

Knowledge Base Creation:

1) Finding 15 Relevant URLs

First, I found a website related to traveling in Japan that covered a wide range of information, from Japan's history and culture to its most popular tourist sites. Most importantly, this website contained external links that would allow me to further traverse the Internet for more information about this topic. The URL of this website is: <https://en.m.wikivoyage.org/wiki/Japan>.

Then, I created a `web_crawler()` function to find fourteen additional relevant links branching off of my initial starter URL. Writing this function required a lot of trial and error in order to ensure that the websites I found contained useful information. To find links, I created a 'link_queue' that initially only contained my starter URL. Then, I scraped the text off of this URL using the Python BeautifulSoup library, including all of its links. I added all of these links to the queue to continue searching indefinitely (although I added a cap of 1000 iterations to prevent infinite loops).

In order for a link to be added to my final 'related_links' list, it had to satisfy the following requirements:

1. The URL could not be in my list of "blocked_hosts".
 - a. After observing the output of my function, I noticed some websites that would not allow me to scrape their text. For example, certain websites included a "Prove you are not a robot" message. To handle these edge cases, I added the website's domain name to a manual list in my program each time I encountered one.
2. It had to have a different hostname than all of the existing links in the array
 - a. This requirement helped ensure diversity in my webpage selection.
3. It had to be successfully scraped by BeautifulSoup.

- a. I used a try/except block to try scraping the website and check whether it was successful.
4. It had to contain at least 500 words of paragraph text (text within <p> tags).
 - a. I specifically searched within the <p> tags using BeautifulSoup to exclude webpage text that was not useful, like text from navigation bars and footers.
 - b. I included the 500 words requirement to ensure that I was scraping enough text to build a sizable knowledge base.
5. The text had to use ASCII characters.
 - a. I added this requirement after I found that a lot of the webpages being returned were written in Japanese.

I continued checking links from my queue using these requirements until I had fifteen valid links.

2) Cleaning the Raw Webpage Text

After extracting the text within the <p> tags using BeautifulSoup, I saved this raw text to its own file in my 'raw_page_texts' directory. Then, I iterated through all of these raw text files to clean them in my clean_text() function. In this function, I clean the text by removing extra whitespace, including new lines and tabs. Then, I used NLTK's sent_tokenize() method to separate the raw text into individual sentences. Then, I iterated through each tokenized sentence and printed it on a line in a new file in my 'clean_page_texts' directory.

3) Determining the Top 25 Terms

To determine the top 25 terms, I applied simple term frequency when analyzing my documents. In my get_top_25_terms() method, I iterated through each cleaned text file from the previous step. Then, I tokenized each individual sentence using the NLTK word_tokenize() method. Finally, I iterated through each token and incremented its count in a dictionary I created.

To ensure my top 25 terms were useful, I applied the following requirements to them:

1. They are alphabetic
2. They are not in the English stop words list from NLTK
3. They have a length greater than five characters

Then, I simply sorted my count dictionary and returned the 25 terms with the greatest counts across all of my documents.

4) Selecting the Top 10 Terms

After finding my top 25 terms, I printed them all to console and manually observed them. I selected ten terms that I thought would be useful for a chatbot travel agent to recognize, and that I thought would result in substantial information in my knowledge base.

Selected terms:

1. japanese
2. country
3. english
4. popular
5. cities
6. hotels
7. restaurants
8. travel
9. international
10. stores

5) Creating the Knowledge Base

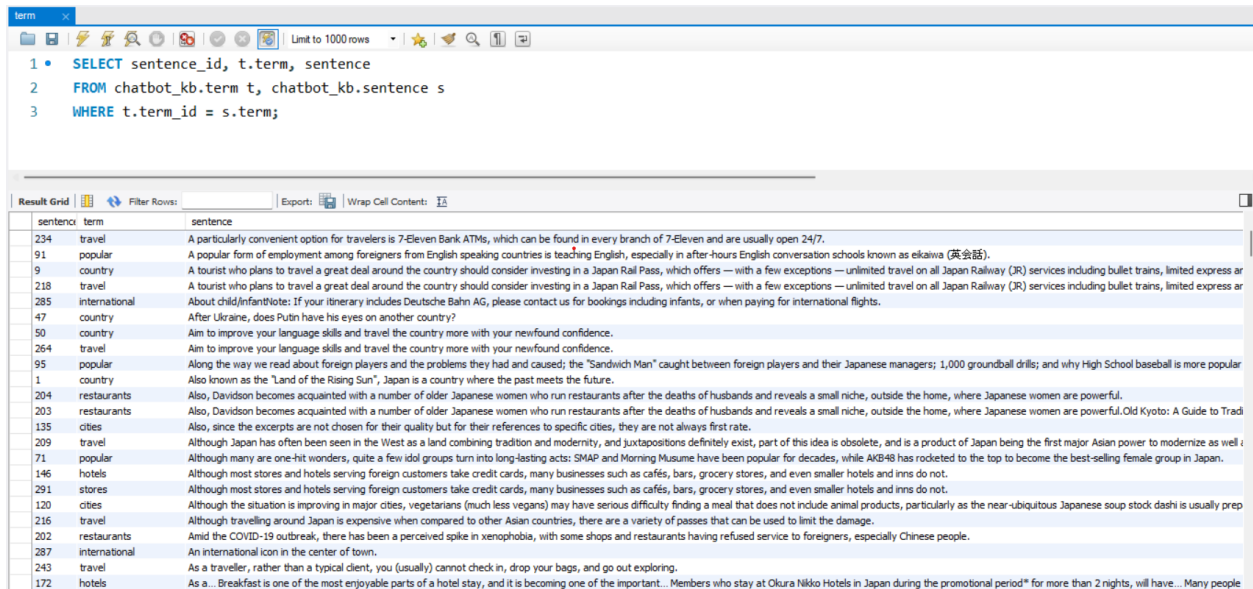
Once I had my top 10 terms, I searched for information related to these terms in my cleaned text files. I created a new dictionary with 10 keys – one for each term listed above. Then, I iterated through each sentence of each cleaned text file. If the sentence contained one of the terms listed above, I included it in the corresponding entry of the dictionary. Then, I dumped this new dictionary into a pickle file named “chatbot_kb.p”.

5) Creating a SQL Database

To further improve the usability of my knowledge base, I created a SQL database for it. The database consists of two tables: term and sentence. The ‘term’ table contains the 10 terms listed above, as well as a unique ID for each. The ‘sentence’ table contains pieces of knowledge that map to these terms. Each entry of the ‘sentence’ table includes a unique ID, a sentence of information (with a maximum length of 1000 characters for handling storage concerns), and a foreign key that maps to the ‘term’ table.

By arranging my data in a SQL database, I can easily query it and view the information. In total, my database contains 310 entries. On the next page, you can find screenshots of my SQL knowledge base showing the execution of a few different queries.

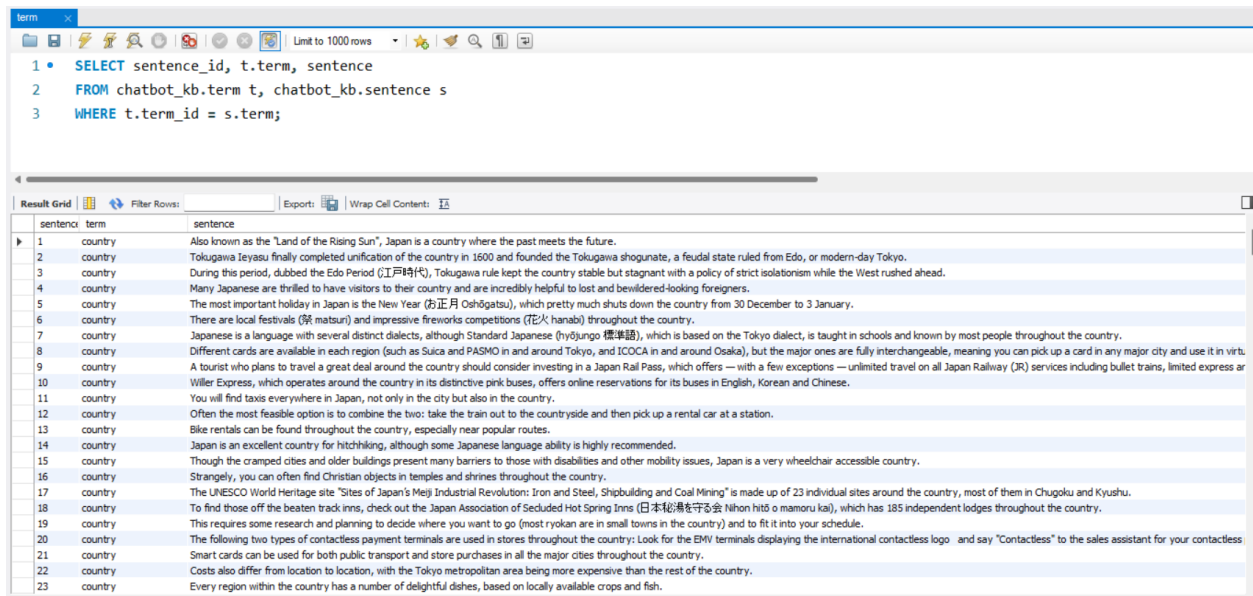
MySQL Database Screenshot, sorted alphabetically:



```
1 • SELECT sentence_id, t.term, sentence
2 FROM chatbot_kb.term t, chatbot_kb.sentence s
3 WHERE t.term_id = s.term;
```

sentence_id	term	sentence
234	travel	A particularly convenient option for travelers is 7-Eleven Bank ATMs, which can be found in every branch of 7-Eleven and are usually open 24/7.
91	popular	A popular form of employment among foreigners from English speaking countries is teaching English, especially in after-hours English conversation schools known as eikaiwa (英会話).
9	country	A tourist who plans to travel a great deal around the country should consider investing in a Japan Rail Pass, which offers — with a few exceptions — unlimited travel on all Japan Railway (JR) services including bullet trains, limited express ar
218	travel	A tourist who plans to travel a great deal around the country should consider investing in a Japan Rail Pass, which offers — with a few exceptions — unlimited travel on all Japan Railway (JR) services including bullet trains, limited express ar
285	international	About child/infant/Note: If your itinerary includes Deutsche Bahn AG, please contact us for bookings including infants, or when paying for international flights.
47	country	After Ukraine, does Putin have his eyes on another country?
50	country	Aim to improve your language skills and travel the country more with your newfound confidence.
264	travel	Aim to improve your language skills and travel the country more with your newfound confidence.
95	popular	Along the way we read about foreign players and the problems they had and caused; the "Sandwich Man" caught between foreign players and their Japanese managers; 1,000 groundball drills; and why High School baseball is more popular
1	country	Also known as the "Land of the Rising Sun", Japan is a country where the past meets the future.
204	restaurants	Also, Davidson becomes acquainted with a number of older Japanese women who run restaurants after the deaths of husbands and reveals a small niche, outside the home, where Japanese women are powerful.
203	restaurants	Also, Davidson becomes acquainted with a number of older Japanese women who run restaurants after the deaths of husbands and reveals a small niche, outside the home, where Japanese women are powerful.Old Kyoto: A Guide to Trad
135	cities	Also, since the excerpts are not chosen for their quality but for their references to specific cities, they are not always first rate.
209	travel	Although Japan has often been seen in the West as a land combining tradition and modernity, and juxtapositions definitely exist, part of this idea is obsolete, and is a product of Japan being the first major Asian power to modernize as well i
71	popular	Although many are one-hit wonders, quite a few idol groups turn into long-lasting acts: SMAP and Morning Musume have been popular for decades, while AKB48 has rocketed to the top to become the best-selling female group in Japan.
146	hotels	Although most stores and hotels serving foreign customers take credit cards, many businesses such as cafés, bars, grocery stores, and even smaller hotels and inns do not.
291	stores	Although most stores and hotels serving foreign customers take credit cards, many businesses such as cafés, bars, grocery stores, and even smaller hotels and inns do not.
120	cities	Although the situation is improving in major cities, vegetarians (much less vegans) may have serious difficulty finding a meal that does not include animal products, particularly as the near-ubiquitous Japanese soup stock dashi is usually prep
216	travel	Although travelling around Japan is expensive when compared to other Asian countries, there are a variety of passes that can be used to limit the damage.
202	restaurants	Amid the COVID-19 outbreak, there has been a perceived spike in xenophobia, with some shops and restaurants having refused service to foreigners, especially Chinese people.
287	international	An international icon in the center of town.
243	travel	As a traveller, rather than a typical client, you (usually) cannot check in, drop your bags, and go out exploring.
172	hotels	As a... Breakfast is one of the most enjoyable parts of a hotel stay, and it is becoming one of the important... Members who stay at Okura Nikko Hotels in Japan during the promotional period* for more than 2 nights, will have... Many people

MySQL Database Screenshot, sorted by term:



```
1 • SELECT sentence_id, t.term, sentence
2 FROM chatbot_kb.term t, chatbot_kb.sentence s
3 WHERE t.term_id = s.term;
```

sentence_id	term	sentence
1	country	Also known as the "Land of the Rising Sun", Japan is a country where the past meets the future.
2	country	Tokugawa Ieyasu finally completed unification of the country in 1600 and founded the Tokugawa shogunate, a feudal state ruled from Edo, or modern-day Tokyo.
3	country	During this period, dubbed the Edo Period (江戸時代), Tokugawa rule kept the country stable but stagnant with a policy of strict isolationism while the West rushed ahead.
4	country	Many Japanese are thrilled to have visitors to their country and are incredibly helpful to lost and bewildered-looking foreigners.
5	country	The most important holiday in Japan is the New Year (お正月 Oshōgatsu), which pretty much shuts down the country from 30 December to 3 January.
6	country	There are local festivals (祭 matsuri) and impressive fireworks competitions (花火 hanabi) throughout the country.
7	country	Japanese is a language with several distinct dialects, although Standard Japanese (hyōjungo 標準語), which is based on the Tokyo dialect, is taught in schools and known by most people throughout the country.
8	country	Different cards are available in each region (such as Suica and PASMO in and around Tokyo, and ICOCA in and around Osaka), but the major ones are fully interchangeable, meaning you can pick up a card in any major city and use it in virtu
9	country	A tourist who plans to travel a great deal around the country should consider investing in a Japan Rail Pass, which offers — with a few exceptions — unlimited travel on all Japan Railway (JR) services including bullet trains, limited express ar
10	country	Willer Express, which operates around the country in its distinctive pink buses, offers online reservations for its buses in English, Korean and Chinese.
11	country	You will find taxis everywhere in Japan, not only in the city but also in the country.
12	country	Often the most feasible option is to combine the two: take the train out to the countryside and then pick up a rental car at a station.
13	country	Bike rentals can be found throughout the country, especially near popular routes.
14	country	Japan is an excellent country for hitchhiking, although some Japanese language ability is highly recommended.
15	country	Though the cramped cities and older buildings present many barriers to those with disabilities and other mobility issues, Japan is a very wheelchair accessible country.
16	country	Strangely, you can often find Christian objects in temples and shrines throughout the country.
17	country	The UNESCO World Heritage site "Sites of Japan's Meiji Industrial Revolution: Iron and Steel, Shipbuilding and Coal Mining" is made up of 23 individual sites around the country, most of them in Chugoku and Kyushu.
18	country	To find those off the beaten track inns, check out the Japan Association of Secluded Hot Spring Inns (日本秘湯を守る会 Nihon hits o mamoru kai), which has 185 independent lodges throughout the country.
19	country	This requires some research and planning to decide where you want to go (most ryokan are in small towns in the country) and to fit it into your schedule.
20	country	The following two types of contactless payment terminals are used in stores throughout the country: Look for the EMV terminals displaying the international contactless logo and say "Contactless" to the sales assistant for your contactless
21	country	Smart cards can be used for both public transport and store purchases in all the major cities throughout the country.
22	country	Costs also differ from location to location, with the Tokyo metropolitan area being more expensive than the rest of the country.
23	country	Every region within the country has a number of delightful dishes, based on locally available crops and fish.

Query for 'hotels':

term		
Limit to 1000 rows		
<pre>1 • SELECT sentence_id, t.term, sentence 2 FROM chatbot_kb.term t, chatbot_kb.sentence s 3 WHERE t.term_id = s.term and t.term = 'hotels';</pre>		
Result Grid		
Filter Rows:	Exports	Wrap Cell Contents: <input type="checkbox"/>
sentence_id	term	sentence
168	hotels	A number of business hotels have Internet access available if you have your own computer, sometimes for free.
146	hotels	Although most stores and hotels serving foreign customers take credit cards, many businesses such as cafés, bars, grocery stores, and even smaller hotels and inns do not.
172	hotels	As a... Breakfast is one of the most enjoyable parts of a hotel stay, and it is becoming one of the important... Members who stay at Okura Nikko Hotels in Japan during the promotional period* for more than 2 nights, will have... Many people
152	hotels	Aside from this, remember that Japan is mostly a cash only country, and credit cards are usually not accepted in smaller forms of accommodation, including small business hotels.
162	hotels	Capsule hotels (ビジネスホテル business hoteu) are usually around ¥10,000 per night and have a convenient location (often near major train stations) as their major selling point, but rooms are usually unbelievably cramped.
155	hotels	Capsule hotels are segregated by sex, and only a few cater to women.
166	hotels	Courier services (宅配便 takuhabin) are useful for sending packages, documents, and even luggage to/from airports, cities, and hotels.
148	hotels	Even bellhops in high end hotels usually do not accept tips.
145	hotels	Even if an accessible room is available, most hotels require booking via phone or email.
161	hotels	Hidden cameras have been found in love hotels, planted by other guests or even occasionally the hotel management.
154	hotels	However, there are several types of uniquely Japanese and far more affordable hotels: Capsule hotels (カプセルホテル kapuseru hoteu) are the ultimate in space-efficient sleeping: for a small fee (normally between ¥3000 and ¥4000), the
150	hotels	In addition to the usual youth hostels and business hotels, you can find several kinds of uniquely Japanese accommodation, ranging from rarefied ryokan inns to strictly functional capsule hotels and utterly over-the-top love hotels.
144	hotels	In major cities, Limousine Buses (リムジンバス rimujin basu) travel from major train stations and hotels to airports.
164	hotels	Local business hotels, farther from major stations, can be significantly cheaper (double room from ¥5000/night).
156	hotels	Many capsule hotels are attached to a spa of varying degrees of luxury and/or legitimacy, often so that entry to the spa costs perhaps ¥2000 but the capsule is only an additional ¥1000.
151	hotels	Many Western hotel booking sites also have only a small selection of Japanese hotels available; to explore the full gamut, use local companies Rakuten Travel or Jalan, which have good English sites.
147	hotels	Many Westernised hotels and restaurants may add a 10% service charge, and family restaurants may add a 10% late-night charge after midnight.
165	hotels	Mimpaku is a great boon for rural areas with few hotels, but in cities the law protects hotels from having too much competition.
174	hotels	okura_hotels Mar 10 nikko_hotels Mar 10 hotel_jalcity Mar 10 nikko_hotels Mar 8 okura_hotels Mar 7 Each of our hotels has its own indi...
153	hotels	On the other hand, three- and four-star business hotels are relatively reasonably priced when compared to prices in major European or North American cities, and even two-star hotels provide impeccable cleanliness and features rarely fou
159	hotels	Once you leave, that is it, so they are not as convenient as proper hotels.
141	hotels	Outside of major tourist attractions and large international hotels, it is rare to find people who are conversant in English.
140	hotels	Peak hanami often coincides with the start of the new school & financial year on April 1, which means lots of people on the move and full hotels in major cities.

Query for 'popular':

term		
Limit to 1000 rows		
<pre>1 • SELECT sentence_id, t.term, sentence 2 FROM chatbot_kb.term t, chatbot_kb.sentence s 3 WHERE t.term_id = s.term and t.term = 'popular';</pre>		
Result Grid		
Filter Rows:	Exports	Wrap Cell Contents: <input type="checkbox"/>
sentence_id	term	sentence
53	popular	The popular board game of Go is also believed to have been introduced to Japan during this period.
54	popular	Chinese influence also reached its peak during the early Heian Period, which saw Buddhism become a popular religion among the masses.
55	popular	Both are popularly elected under a parallel system, where some seats are filled by individual candidates and others are filled by a party.
56	popular	Low-cost carriers have become increasingly popular with cheap domestic and international flights, with companies such as Jetstar (Australia), Skymark, and Peach (Osaka) offering competition to JAL and ANA.
57	popular	Bike rentals can be found throughout the country, especially near popular routes.
58	popular	Hiking is very traditional and popular in Japan.
59	popular	On the most popular route, you will need to use your hands for support, but no actual climbing is required; you can easily climb Fuji with just adequate clothing, some basic gear (sunscreen, headlamp, etc.
60	popular	Golf is popular with the Japanese.
61	popular	Surfing is somewhat popular, as the surf can be very good on both coasts (during typhoon season [Aug-Oct] on the Pacific coast, and during winter on the Sea of Japan coast).
62	popular	Baseball (野球 yakyū) has been hugely popular ever since it was introduced to Japan in the 1870s by an American professor.
63	popular	Baseball fans travelling internationally may find Japan to be one of the great examples of baseball popularity outside of the United States.
64	popular	Tickets to baseball games are generally easy to get, even on the day of a game, although popular games should of course be reserved in advance.
65	popular	Soccer (サッカー sakka) is also popular in Japan, although it plays second fiddle to baseball.
66	popular	Sumo wrestling (相撲 sumō) is a popular Japanese sport.
67	popular	Professional wrestling (プロレス puroresu) also enjoys major popularity.
68	popular	Besides Go, another popular board game in Japan is shogi (将棋 shōgi) or Japanese chess.
69	popular	Western classical music (クラシック[音楽] kurashiku [ongaku]) is moderately popular in Japan with people of all ages.
70	popular	Jazz (ジャズ jazz) has been very popular in Japan since the 1930s.
71	popular	Although many are one-hit wonders, quite a few idol groups turn into long-lasting acts: SMAP and Morning Musume have been popular for decades, while AKB48 has rocketed to the top to become the best-selling female group in Japan.
72	popular	Music festivals (ロックフェスティバル) rokku fesutibaru, shortened to ロックフェス rokku fesu or just フェス, fesu) are also popular, drawing tens of thousands of people.
73	popular	You may be more interested in the less well-known taishū engaki (pop theater) or modern comedy, such as rakugo solo storytellers, extremely popular manzai stand-up duos, or Western-style comedy.
74	popular	Kabuki (歌舞伎) is a popular type of dance-drama.
75	popular	Much less well-known is taishū engaki (大衆演劇), a vague term meaning "theater for the masses" or "popular theater".

Query for 'restaurants':

```

1 • SELECT sentence_id, t.term, sentence
2   FROM chatbot_kb.term t, chatbot_kb.sentence s
3  WHERE t.term_id = s.term and t.term = 'restaurants';

```

Result Grid

Filter Rows

Export

Wrap Cell Content

	sentence_id	term	sentence
▶	175	restaurants	Christmas Eve is considered to be one of the most romantic days of the year in Japan, and restaurants will be fully booked by young couples looking to have a romantic night out, so be sure to make your dinner reservations well in advance.
	176	restaurants	In addition to public transport, smart cards are used for all sorts of electronic payments, so they can be used at vending machines, convenience stores, fast food restaurants, etc.
	177	restaurants	Hostesses work in bars and sing karaoke to entertain, compared to geisha coming to tea houses and restaurants to perform traditional Japanese arts.
	178	restaurants	Maid cafés and other cosplay restaurants have employees dressed as French maids pamper their clients while serving them beverages and food.
	179	restaurants	This includes vending machines, convenience stores, fast food restaurants, etc.
	180	restaurants	Many Westernised hotels and restaurants may add a 10% service charge, and family restaurants may add a 10% late-night charge after midnight.
	181	restaurants	The Michelin Guide is considered by many Western visitors to be the benchmark of good restaurants in Japan.
	182	restaurants	But many top fine dining restaurants are not listed in it by choice.
	183	restaurants	Many restaurants give you a hot towel (o-shibori) to wipe your hands with (not your face) as soon as you sit down.
	184	restaurants	In all types of Japanese restaurants, staff generally ignore you until you ask for something.
	185	restaurants	Tipping is not customary in Japan, although many sit-down restaurants apply 10% service charges and 24-hour "family restaurants" usually have a 10% late-night surcharge.
	186	restaurants	The number of restaurants (レストラン resutoran) in Japan is stupendous, and you will never run out of places to go.
	187	restaurants	Japan is tied with France for first place as the country with the most Michelin-star restaurants.
	188	restaurants	Most sushi restaurants keep a handy multilingual decoding key on hand or on the wall.
	189	restaurants	At the finest sushi restaurants, the chef puts a dab of fiery wasabi radish into the sushi, and glazes the fish with soy sauce for you.
	190	restaurants	Most restaurants, though, provide them at the table.
	191	restaurants	Tempura has entered the Japanese fine dining repertoire, and there are numerous fine tempura omakase restaurants in which the chef deep-fries the dish in front of you and puts it directly on your plate to be eaten immediately.
	192	restaurants	Curry restaurants can be found throughout the country, but there are also localized versions, such as Kanazawa Curry, Bizen Curry, and Kuwana Curry.
	193	restaurants	You can also find restaurants in department stores, often on the top floors.
	194	restaurants	Japanese fast food restaurants offer decent quality at reasonable prices.
	195	restaurants	There are also a number of Japanese family restaurants (ファミリーレストラン famiri resutoran), serving a wide variety of dishes, including steak, pasta, Chinese style dishes, sandwiches, and other foods.

Query for 'hotels':

term

6) Sample Dialog

Using pieces of knowledge I found in my database, here is what I envision could be a chatbot conversation. All of the Chatbot responses come directly from my knowledge base with no modification of the text. I also highlighted the keywords that the user said that would help the chatbot form its answer using my knowledge base.

Chatbot: Hello, how can I help you today?

User: I want to plan a trip to Japan! What is the best way to **travel** there?

Chatbot: Delta Air Lines, United Airlines, and American Airlines also operate flights from numerous U.S. cities into Narita or Haneda, as does Air Canada from several Canadian cities.

User: What is the best way to **travel** around Japan?

Chatbot: A tourist who plans to travel a great deal around the country should consider investing in a Japan Rail Pass, which offers — with a few exceptions — unlimited travel on all Japan Railway (JR) services including bullet trains, limited express and regular commuter trains.

User: What are the most **popular** things to do in Japan?

Chatbot: Hiking is very traditional and popular in Japan.

User: Hiking sounds fun! Where should I hike?

Chatbot: On the most popular route, you will need to use your hands for support, but no actual climbing is required; you can easily climb Fuji with just adequate clothing, some basic gear (sunscreen, headlamp, etc.

User: What about food? Does Japan have good **restaurants**?

Chatbot: Japan is tied with France for first place as the country with the most Michelin-star restaurants.

User: That sounds great! I want to go to Japan now. Can you help me book a **hotel**?

Chatbot: Some major chains of cheaper business hotels include Tokyu REI Hotels, known for its generously sized rooms, Sunroute Hotels and Toyoko Inn. Even if an accessible room is available, most hotels require booking via phone or email.