

```
1 import re
2 from importlib.resources.readers import remove_duplicates
3
4
5 class Data:
6
7     # list of possible type of questions ~ truncated for flexibility [non-exhaustive]
8     possibleList = [{"capital", "distanc", "weather", "movi", "forecast", "cit", "length", "climat", "humidit", "director", "actor"}, # direct-answer question
9                     ["exercis", "diet", "cook", "workout", "routin", "gym", "activit", "nutri", "wellness", "recipi", "fitnes", "yoga", "meditat", "stretch", "cardio", "strength", "vitamin", "calori"], # health-related questions
10                    ["calendar", "remind", "task", "schedul", "event", "deadlin", "project", "checklist", "alert", "notif", "organ", "priorit", "goal", "plann", "timelin", "focus", "track", "habit", "workflow"], # productivity questions
11                    ["scor", "gam", "jok", "song", "challeng", "puzzl", "music", "lyric", "match", "adventur", "humor", "quiz", "fun", "comed", "story", "celebr", "sport", "trend"], # entertainment questions
12                    ["plu", "minus", "multipl", "divid", "formula", "concept", "ratio", "algebra", "geometr", "calculus", "integrat", "deriv", "vector", "probabil", "statist", "measur", "equation", "matrix", "quantit"], # mathematical questions
13                    ["pric", "mean", "fact", "happen", "latest", "explain", "differenc", "orig", "reason", "impact", "histor", "overview", "background"], # knowledge-building question
14                    ["best", "advic", "help", "tip", "plan", "stuck", "assist", "recommend", "suggest", "guid", "strategy", "solv", "improv", "overcom", "choic", "option"]] # advice-seeking questions
15
16     # list of possible types of questions ~ original unmodified [non-exhaustive]
17     dictionaryList = [{"capital", "distance", "weather", "movie", "forecast", "city", "length", "climate", "humidity", "director", "actor"},
18                      ["exercise", "diet", "cook", "workout", "routine", "gym", "activity", "nutri", "wellness", "recipie", "fitness", "yoga", "meditate", "stretch", "cardio", "strength", "vitamin", "calories"],
19                      ["calendar", "remind", "task", "schedule", "event", "deadline", "project", "checklist", "alert", "notif", "organ", "prioritize", "goal", "plann", "timeline", "focus", "track", "habit", "workflow"],
20                      ["score", "game", "joke", "song", "challenge", "puzzle", "music", "lyric", "match", "adventure", "humor", "quiz", "fun", "comedian", "story", "celebrate", "sport", "trend"],
21                      ["plu", "minus", "multiple", "divide", "formula", "concept", "ratio", "algebra", "geometry", "calculus", "integrate", "derivative", "vector", "probability", "statistic", "measure", "equation", "matrix", "quantity"],
22                      ["price", "mean", "fact", "happen", "latest", "explain", "difference", "origin", "reason", "impact", "history", "overview", "background"],
23                      ["best", "advice", "help", "tip", "plan", "stuck", "assist", "recommend", "suggest", "guide", "strategy", "solve", "improve", "overcome", "choice", "option"]]
24
25     # list of possible cities that the user might reference [non-exhaustive]
26     specificPlaceList = ["Los Angeles", "Chicago", "San Francisco", "Miami", "Austin", "Las Vegas", "Paris", "London", "Tokyo", "Sydney", "Rome", "Barcelona", "Berlin", "Dubai", "Toronto", "Seoul", "Bangkok", "Mexico City", "Riverside", "Cape Town", "California", "Florida", "Texas", "New York", "Nevada", "Hawaii", "Colorado", "Alaska", "Arizona", "Utah", "Illinois", "Michigan", "Washington", "Georgia", "North Carolina", "Tennessee", "South Carolina", "Oregon", "New Jersey", "Virginia", "United States", "Canada", "United Kingdom", "France", "Italy", "Spain", "Mexico", "Germany", "Australia", "Brazil", "Japan", "India", "South Korea", "Thailand", "South Africa", "China", "Russia", "Egypt", "Argentina", "New Zealand", "Pakistan"]
27
28
29
30
31     # list of possible pop-culture references [non-exhaustive]
32     specificPopCultureList = [{"The Avengers", "Star Wars", "The Matrix", "Harry Potter", "Jurassic Park", "Titanic", "The Godfather", "Pulp Fiction", "Back to the Future", "The Lion King", "Apple", "Nike", "Tesla", "Coca-Cola", "McDonald's", "Amazon", "Google", "Adidas", "Disney", "Microsoft", "Friends", "Game of Thrones", "The Office", "Stranger Things", "Breaking Bad", "The Simpsons", "The Mandalorian", "The Crown", "The Walking Dead", "Westworld", "The Beatles", "Beyonce", "Kanye West", "Taylor Swift", "Elvis Presley", "Michael Jackson", "Ariana Grande", "Drake", "Lady Gaga", "Eminem", "Super Mario Bros.", "Minecraft", "Fortnite", "The Legend of Zelda", "Call of Duty", "Grand Theft Auto", "Pokémon", "League of Legends", "FIFA", "The Witcher", "Rolls-Royce", "Ferrari", "Lamborghini", "Porsche", "Maserati", "Bentley", "Aston Martin", "Bugatti", "McLaren", "Mercedes-Benz", "Lexus", "BMW", "Audi"}]
33
34
35
36
37
38
39     # word that is autocorrected
40     correctedWord = None
41
42     # method that removes accidental duplicates found by regular expression
43     def __remove_duplicate(self, list):
44         accList = []
45         for i in list:
46             if i not in accList:
47                 accList.append(i)
48         return accList
49
50
51     # uses Python Regular Expressions to derive key data in a structural format, replicating a basic version of Natural Language Processing
52     # "QT" = Question Type && "I" = Identifier
53     def parsedData(self, userInput, original_input):
54         question_type = self.__remove_duplicate(re.findall(r"(what|who|why|where|when|how|will|can|play|lets|let|should|is|tell|give|if)", userInput))
55         identifiers = self.__remove_duplicate(re.findall(r"([capital|best|cit|length|climat|humidit|director|actor|task|schedul|event|deadlin|project|checklist|alert|notif|organ|advic|stuck|help|tip|distanc|plan|weather|forecast|latest|r]|happen|movi|exercis|song|diet|workout|explain|differenc|routin|gym|activit|nutri|wellness|recipi|fitnes|calendar|remind|cook|scor|pric|mean|plu|ratio|minus|multipl|divid|r]|jok|gam|fact|formula|concept|algebra|geometr|challeng|puzzl|music|lyric|match|adventur|humor|yoga|meditat|stretch|cardio|r]|strength|vitamin|calori|priorit|goal|plann|timelin|focus|track|habit|workflow|quiz|fun|comedi|story|celebr|sport|trend|r]|calculus|integrat|deriv|vector|probabil|statist|measur|equation|matrix|quantit|orig|reason|impact|histor|overview|background|assist|recommend|suggest|guid|strategy|solv|improv|overcom|choic|option)", userInput))
56
57         print(question_type)
58         print(identifiers)
59         # returns error if either one of the variables above is empty, else, normal dictionary returned
60         if len(question_type) < 1 or len(identifiers) < 1:
61             # extra step of autocorrect for better interpretation
62             issue = {"QT": question_type, "I": self.autoCorrect(userInput, self.possibleList).split()}
63             if self.autoCorrect(userInput, self.possibleList) != "ERROR":
64                 self.correctedWord = self.autoCorrect(original_input, self.dictionaryList)
65                 print("I am assuming that you mean to say " + self.autoCorrect(original_input, self.dictionaryList) + ". The original word is autocorrected.\n")
66             return issue
67         else:
68             result = {"QT": question_type, "I": identifiers}
69             self.correctedWord = identifiers[0]
70             return result
71
72     # takes the original sentence inputted by the user and then removes suffixes; local change not global
73     def stemWord(self, userInput):
74         return re.sub(r'\b(?:\b|)(e|es|ingled|s|se|ication|ization|isation|ized|ised|ied|ous|y|ies|tion|ent|ents|er|ers|ic|ation|ating|ize|ian|ate|ative|atives|ity|ics|in|inate|ance)\b', "", userInput)
75
76     # finds any key data similar to the list above that might be misspelled to reinterpret input ~ similarity >= 70%
77     def autoCorrect(self, userInput, c_list):
78         inputList = userInput.split()
79         count = 0
80         for r in range(len(c_list)):
81             for c in range(len(c_list[r])):
82                 for n in inputList:
83                     iter = min(len(n), len(c_list[r][c]))
84                     for i in range(iter):
85                         if list(n).__getitem__(i) == list(c_list[r][c]).__getitem__(i):
86                             count+=1
87                         elif (i + 1) < len(c_list[r][c]):
88                             if list(n).__getitem__(i) == list(c_list[r][c]).__getitem__(i + 1):
89                                 count+=1
90         if (count / max(len(n), len(c_list[r][c]))) * 100 >= 70:
91             return c_list[r][c]
92         count = 0
```

```
96         return "ERROR"
```