

1 Exercise 1

1.1 1a & 1b

This Bloom Filter spell correction system is initialized by inserting a large dictionary of valid words using multiple hash functions. When a misspelled word is encountered, the system generates possible corrections by simulating all single-character substitutions. Each candidate word is checked against the Bloom filter, and if it is flagged as a possible match, it is added to the list of suggested corrections. Then we evaluate the effectiveness of these suggestions based on two key criteria: the suggestion list should contain no more than three words and must include the correct spelling. Also, the other evaluation metric we need to calculate is called misidentified percentage, which is defined as the percentage of the cases where typed words are in the bloom filter while not correct words.

```
1 from bitarray import bitarray
2 from hashlib import sha3_256, sha256, blake2b
3 import json
4 import string
5
6 # Bloom Filter
7 class BloomFilter:
8     def __init__(self, size, num_hashes):
9         self.size = size
10        self.num_hashes = num_hashes
11        self.bit_array = bitarray(size)
12        self.bit_array.setall(0)
13        self.hash_functions = [self.my_hash, self.my_hash2, self.my_hash3] #three hashes
14
15    def my_hash(self, s):
16        return int(sha256(s.lower().encode()).hexdigest(), 16) % self.size
17
18    def my_hash2(self, s):
19        return int(blake2b(s.lower().encode()).hexdigest(), 16) % self.size
20
21    def my_hash3(self, s):
22        return int(sha3_256(s.lower().encode()).hexdigest(), 16) % self.size
23
24    def add(self, word):
25        for i in range(self.num_hashes):
26            index = self.hash_functions[i](word)
27            self.bit_array[index] = True
28
29    def might_contain(self, word):
30        for i in range(self.num_hashes):
31            index = self.hash_functions[i](word)
32            if not self.bit_array[index]:
33                return False
34        return True
35
36 # spell check using single-character substitutions
37 def spell_check(word, bloom_filter):
38     alphabet = string.ascii_lowercase
39     suggestions = set()
40
41     # Single character substitutions, we regard these substitutions as potential suggestions
42     # , we will check if they are real suggestions
43     for i in range(len(word)):
44         for char in alphabet:
45             new_word = word[:i] + char + word[i+1:]
46             if bloom_filter.might_contain(new_word):
47                 suggestions.add(new_word)
48
49     return list(suggestions)
50
51 # Evaluate results
52 def evaluate_performance(bloom_filter, dataset):
53     typo_total = 0
54     good_suggestions = 0
```

```

54 misidentified_count = 0
55
56 for typed_word, correct_word in dataset:
57     if typed_word != correct_word: # the typed word is a misspelled one if they are not
        identical.
58         typo_total += 1
59         suggestions = spell_check(typed_word, bloom_filter)
60
61         # Good suggestion
62         if len(suggestions) <= 3 and correct_word in suggestions:
63             good_suggestions += 1
64
65         # Misidentified
66         if bloom_filter.might_contain(typed_word):
67             misidentified_count += 1
68
69     # Good suggestion
70     good_suggestion_percentage = (good_suggestions / typo_total) * 100
71
72     # Misidentified percentage: cases where typed_word is in the Bloom filter but is not
        correct_word
73     misidentified_rate = (misidentified_count / typo_total) * 100
74
75     return good_suggestion_percentage, misidentified_rate
76
77 size = int(1e7)
78 num_hashes = 3
79 bloom_filter = BloomFilter(size, num_hashes)
80
81 with open('words.txt') as f:
82     for line in f:
83         word = line.strip()
84         bloom_filter.add(word)
85 #Use the given example for experiment
86 typo = "floeer"
87 suggestions = spell_check(typo, bloom_filter)
88 print(f"Suggestions for '{typo}': {suggestions}")
89
90 with open('typos.json') as f:
91     dataset = json.load(f)
92
93 performance = evaluate_performance(bloom_filter, dataset)
94 print(f"Performance: {performance}%")

```

Listing 1: Bloom Filter for spell check

This 1 is the result when we set the size as 10^7 and number of hash functions =3.

```

Suggestions for 'floeer': ['flower', 'floter']
Performance: (91.74799999999999, 0.20400000000000001)%

```

Figure 1: 1b Experiment result

1.2 1c

Then we can plot the results using the code below.

```

1 import matplotlib.pyplot as plt
2
3 def experiment_and_plot(bloom_filter_sizes, num_hash_functions_list, dataset):
4     results = []
5
6     # Experiment with varying filter sizes and number of hash functions
7     for size in bloom_filter_sizes:
8         for num_hashes in num_hash_functions_list:
9             # Initialize a new Bloom filter with the specified size and number of hash
                functions
10             bloom_filter = BloomFilter(size, num_hashes)
11
12             with open('words.txt') as f:
13                 for line in f:
14                     word = line.strip()
15                     bloom_filter.add(word)
16

```

```

17     # Evaluate the performance
18     good_suggestion_percentage, misidentified_rate = evaluate_performance(
        bloom_filter, dataset)
19
20     # Store results
21     results.append({
22         'size': size,
23         'hashes': num_hashes,
24         'good_suggestions': good_suggestion_percentage,
25         'misidentified': misidentified_rate
26     })
27     plot_combined_results(results)
28
29 ## Single y-axis for both metrics
30 def plot_combined_results(results):
31     sizes = sorted(list(set([result['size'] for result in results])))
32     num_hashes_list = sorted(list(set([result['hashes'] for result in results])))
33
34     plt.figure()
35
36     # Plot both Good Suggestions and Misidentified on the same axis
37     #I ask ChatGPT for help and it separates the good_suggestion and misidentified from the
    results list.
38     for num_hashes in num_hashes_list:
39         x = [result['size'] for result in results if result['hashes'] == num_hashes]
40         y_good = [result['good_suggestions'] for result in results if result['hashes'] ==
        num_hashes]
41         y_misidentified = [result['misidentified'] for result in results if result['hashes']
        == num_hashes]
42
43         plt.plot(x, y_good, marker='o', linestyle='--', label=f'{num_hashes} Hash Functions -
        Good Suggestions')
44         plt.plot(x, y_misidentified, marker='x', linestyle='--', label=f'{num_hashes} Hash
        Functions - Misidentified')
45
46     plt.xscale('log')
47     plt.xlabel('Bloom Filter Size (log scale)')
48     plt.ylabel('Percentage')
49     plt.title('Effect of Bloom Filter Size and Hash Functions on Good Suggestions and False
        Positives')
50     plt.grid(True)
51     plt.legend(loc='center left', bbox_to_anchor=(0, 0.5))# arrage the position of the
    legend
52
53     plt.tight_layout()
54     plt.show()
55
56 # Experiment parameters
57 bloom_filter_sizes = [int(1e1),int(1e2),int(1e3),int(1e4),int(1e5), int(1e6), int(1e7), int
    (1e8), int(1e9), int(1e10)] # Filter sizes
58 num_hash_functions_list = [1, 2, 3] # Number of hash functions
59
60 # Load the dataset (typos.json)
61 with open('typos.json') as f:
62     dataset = json.load(f)
63 # Run the experiment and plot results
64 experiment_and_plot(bloom_filter_sizes, num_hash_functions_list, dataset)

```

Listing 2: plot results of different sizes and number of hash functions

The final result of the implementation of our bloom filter is shown in 2.

1.3 1d

From the figure, we can find that the bloom filter with 1 hash function with a size larger than 10^9 can reach the 85% good suggestion rate. The bloom filter with 2 hash functions with a size larger than 10^8 can reach the 85% good suggestion rate. The bloom filter with 3 hash functions with a size larger than 10^7 can reach the 85% good suggestion rate.

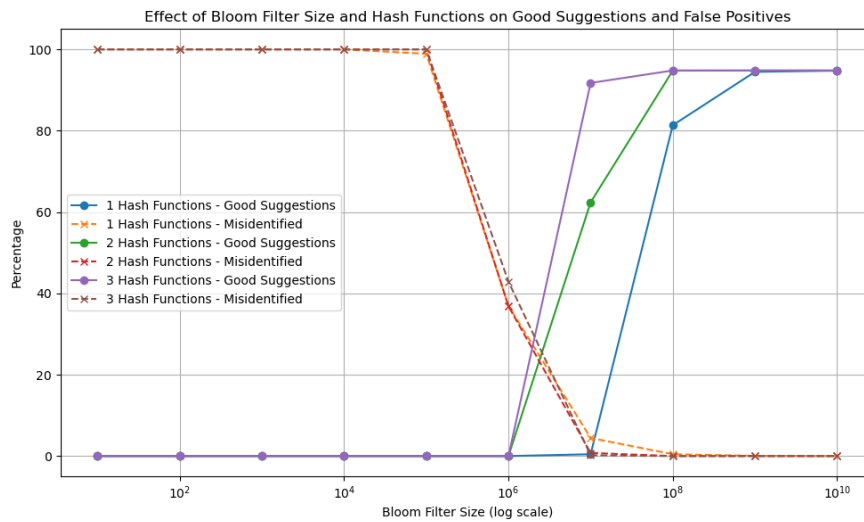


Figure 2: 1c Output

2 Exercise 2

2.1 2a

To make our merge sort algorithm capable of dictionary type data, we add a new parameter called key to search each key of the dictionary iteratively.

```

1 def alg2_key(data, key):
2     if len(data) <= 1:
3         return data
4     else:
5         split = len(data) // 2
6         left = iter(alg2_key(data[:split], key))
7         right = iter(alg2_key(data[split:], key))
8         result = []
9
10        left_top = next(left)
11        right_top = next(right)
12
13        while True:
14            if left_top[key] < right_top[key]:
15                result.append(left_top)
16                try:
17                    left_top = next(left)
18                except StopIteration:
19                    # nothing remains on the left; add the right + return
20                    return result + [right_top] + list(right)
21            else:
22                result.append(right_top)
23                try:
24                    right_top = next(right)
25                except StopIteration:
26                    # nothing remains on the right; add the left + return
27                    return result + [left_top] + list(left)

```

Listing 3: Modification of merge sort

As shown in figure 3, we first build up a simple dictionary containing the data, and the new merge sort algorithm successfully sorts the data based on the patient ID.

```

data = [
    {'patient_id': 3, 'patient_data': 'A'},
    {'patient_id': 1, 'patient_data': 'B'},
    {'patient_id': 4, 'patient_data': 'C'},
    {'patient_id': 2, 'patient_data': 'D'}
]

sorted_data = alg2_key(data, 'patient_id')
print(sorted_data)

[{'patient_id': 1, 'patient_data': 'B'}, {'patient_id': 2, 'patient_data': 'D'}, {'patient_id': 3, 'patient_data': 'A'}, {'patient_id': 4, 'patient_data': 'C'}]

```

Figure 3: 2a Output

2.2 2b

First of all, we should generate our data. We randomly choose an integer as the patient id and randomly generate a string as the patient data. It is likely to have duplicates, but the probability is low and it does not influence the performance of sorting.

```
1 def generateNumbers(numElements):
2     randomArray = []
3     for i in range(numElements):
4         patient_id = random.randint(1, 1000000000)
5         patient_data = ''.join(random.choices(string.ascii_letters + string.digits, k=4))
6         randomArray.append({'patient_id': patient_id, 'patient_data': patient_data})
7     return randomArray
```

Listing 4: Generate Data

Here is our entire implementation of this task. After generating data, we designed a parallel sorting. The parallel sort function performs a merge sort on an array using parallel processing. Also, if the array size is smaller than a given 'threshold', it switches to a serial sort. For larger arrays, it splits the array into two halves, sorts them in parallel using multiple processes. The *apply_async* method is used to asynchronously run the sorting on both halves, allowing them to execute concurrently. Once both halves are sorted, they are merged back together using a standard merge step (alg2 key), which is a serial merge algorithm.

```
1 from time import perf_counter
2 import random
3 import multiprocessing as mp
4 import matplotlib.pyplot as plt
5 import numpy as np
6 import string
7
8 # Logging performance data to a file
9 def logFile(seqtime1, seqtime2, paratime1, paratime2, speedup):
10     with open("merge.log", "a") as f:
11         f.write("Sequential: \n")
12         f.write("Elapsed time: ")
13         f.write('{:.20f}'.format(seqtime1 - seqtime2) + "\n")
14         f.write("Parallel: \nElapsed time: ")
15         f.write('{:.20f}'.format(paratime1 - paratime2) + "\n")
16         f.write("Speedup: ")
17         f.write('{:.10f}'.format(speedup) + "\n")
18
19 # Start Timer
20 def startTime():
21     return perf_counter()
22
23 # End Timer
24 def endTime():
25     return perf_counter()
26
27 # Custom merge sort for dictionary data based on 'patient_id'
28 def alg2_key(data, key):
29     if len(data) <= 1:
30         return data
31     else:
32         split = len(data) // 2
33         left = iter(alg2_key(data[:split], key))
34         right = iter(alg2_key(data[split:], key))
35         result = []
36
37         left_top = next(left)
38         right_top = next(right)
39
40         while True:
41             if left_top[key] < right_top[key]:
42                 result.append(left_top)
43                 try:
44                     left_top = next(left)
45                 except StopIteration:
46                     return result + [right_top] + list(right)
47             else:
48                 result.append(right_top)
49                 try:
50                     right_top = next(right)
51                 except StopIteration:
52                     return result + [left_top] + list(left)
53
```

```

54 # Generates randomly a number of elements in dictionary form with 'patient_id' and '
    patient_data'
55 def generateNumbers(numElements):
56     randomArray = []
57     for i in range(numElements):
58         patient_id = random.randint(1, 100000000)
59         patient_data = ''.join(random.choices(string.ascii_letters + string.digits, k=4))
60         randomArray.append({'patient_id': patient_id, 'patient_data': patient_data})
61     return randomArray
62
63 def copyArray(arr):
64     return arr.copy()
65
66 def sumArray(arr):
67     return sum(arr)
68
69 # Parallel sorting function (using multiprocessing) for dictionary data
70 def parallel_sort(arr, key, threshold=1000):
71     if len(arr) <= threshold:
72         return alg2_key(arr, key) # Switch to serial sort for small arrays
73
74     split = len(arr) // 2
75     with mp.Pool(processes=2) as pool: # Using pool for parallel processing
76         left_future = pool.apply_async(parallel_sort, (arr[:split], key))
77         right_future = pool.apply_async(parallel_sort, (arr[split:], key))
78         left = left_future.get()
79         right = right_future.get()
80
81     return alg2_key(left + right, key)
82
83 # Plot performance
84 def plot_performance(sizes, serial_times, parallel_times):
85     plt.figure(figsize=(10, 6))
86     plt.loglog(sizes, serial_times, label='Serial Merge Sort', marker='o')
87     plt.loglog(sizes, parallel_times, label='Parallel Merge Sort', marker='o')
88     plt.xlabel('Dataset Size (log scale)')
89     plt.ylabel('Time (log scale)')
90     plt.title('Performance of Serial vs Parallel Merge Sort')
91     plt.legend()
92     plt.grid(True)
93     plt.show()
94
95 if __name__ == '__main__':
96     sizes = [10**i for i in range(5,8)] # Test sizes from 10^5 - 10^7
97     sTimeArr = []
98     pTimeArr = []
99
100     for size in sizes:
101         arr1 = generateNumbers(size)
102         arr2 = copyArray(arr1)
103
104         # Sequential merge sort
105         start = startTime()
106         alg2_key(arr1, 'patient_id') # Perform serial sort based on 'patient_id'
107         end = endTime()
108         seq_time = end - start
109         sTimeArr.append(seq_time)
110         print(f"Sequential Sort - Size {size}: {seq_time:.6f} seconds")
111
112         # Parallel merge sort
113         pStart = startTime()
114         p = mp.Process(target=parallel_sort, args=(arr2, 'patient_id'))
115         p.start()
116         p.join()
117         pEnd = endTime()
118         par_time = pEnd - pStart
119         pTimeArr.append(par_time)
120         print(f"Parallel Sort - Size {size}: {par_time:.6f} seconds")
121
122         # Check for 70% time reduction
123         if par_time <= 0.7 * seq_time:
124             print(f"Parallel run for size {size} met the 70% criterion.")
125         if par_time <= 0.33 * seq_time:
126             print(f"Parallel run for size {size} met the 1/3 criterion for extra credit.")
127
128     # Plot the results

```

Listing 5: Parallel Merge Sort

The result is shown in figure 4. As we can see, serial merge sort uses less time with a small data size. When the data size increased exponentially, the parallel merge sort consume time much less than the serial merge sort. To be more specific, when the dataset size is 10^5 , the time consumed by serial merge sort is 0.122866, and the time consumed by parallel merge sort is 0.685740; When the size is 10^6 , the time consumed by serial merge sort is 1.877798, the time consumed by parallel merge sort is 1.212943 (**less than 70% time of the serial merge sort**); When the size is 10^7 , the time consumed by serial merge sort is 30.003212, the time consumed by parallel merge sort is 9.570690 (**less than 1/3 time of the serial merge sort**)

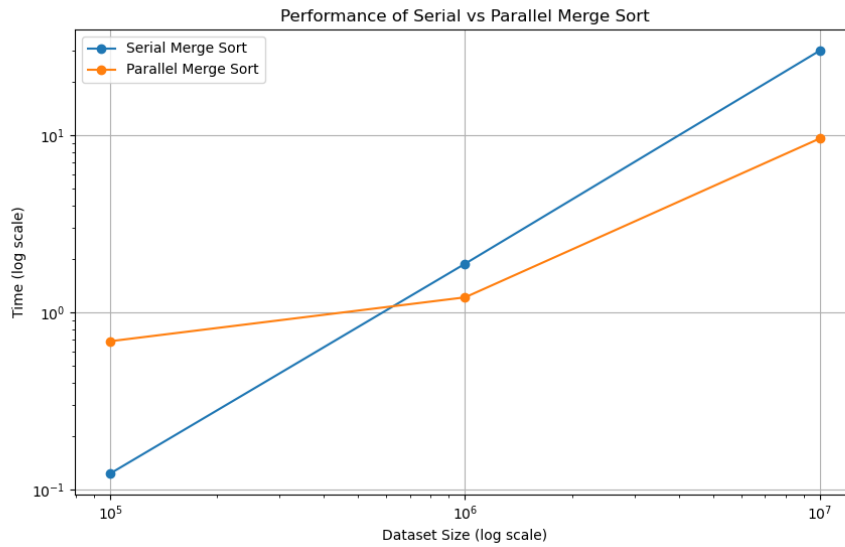


Figure 4: 2b Output

3 Exercise 3

3.1 3a

The 'get id' function queries the PubMed API for research articles related to a specific disease published in 2023. It makes a GET request to retrieve up to 1000 PubMed IDs, parses the XML response, and returns a list of these IDs. If the request fails, it prints an error message and returns an empty list.

```
1 import requests
2 import xml.etree.ElementTree as ET
3 import time
4
5 def get_id(disease):
6     # Make the request to PubMed API
7     r = requests.get(
8         "https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?"
9         f"db=pubmed&term={disease}+AND+2023[pdat]&retmode=xml&retmax=1000"
10    )
11    # Respect rate limit as suggested
12    time.sleep(1)
13
14    # Check for successful response
15    if r.status_code == 200:
16        # Parse the XML response
17        root = ET.fromstring(r.content)
18
19        # Extract all PubMed IDs
20        pubmed_ids = [id_elem.text for id_elem in root.findall(".//Id")]
21        return pubmed_ids
22    else:
23        # Error handling
24        print("Failed to fetch IDs:", r.status_code)
25        return []
```

Listing 6: get id

3.2 3b, 3c & 3d

The 'get metadata' function retrieves metadata (article title and abstract) for a list of PubMed IDs by making a POST request to the PubMed API. It sends the list of IDs, receives the metadata in XML format, and parses it using the 'ElementTree' module. For each article, the function extracts the title and abstract (including handling subsections of abstracts), then stores this information in a dictionary, associating it with the corresponding PubMed ID and the original query term. The helper function 'get text' is used to extract text content from the XML nodes. If no IDs are provided, the function returns an empty dictionary.

```
1 def get_text(element):
2     """Helper function to retrieve text content from an ElementTree node and its children.
3     """
4     return ''.join(element.itertext())
5
6 def get_metadata(id_list, query_term):
7     metadata_dict = {}
8
9     if id_list:
10        # Make the POST request to fetch metadata
11        r = requests.post(
12            "https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi",
13            data={
14                "db": "pubmed",
15                "retmode": "xml",
16                "id": ",".join(id_list)
17            }
18        )
19
20        # Respect API rate limiting
21        time.sleep(1)
22
23        if r.status_code == 200:
24            # Parse the XML response using ElementTree and iterate through PubmedArticle
25            # nodes
26            root = ET.fromstring(r.content)
27            articles = root.findall(".//PubmedArticle")
28
29            for article, pubmed_id in zip(articles, id_list):
```



```

28         # Extract title using ElementTree
29         title_elem = article.find(".//ArticleTitle")
30         title = get_text(title_elem) if title_elem is not None else "" #Some
abstract context is missing
31
32         # Extract abstract using ElementTree, and handle multiple parts of
structured abstracts
33         abstract_elems = article.findall(".//AbstractText")
34         abstract = " ".join(get_text(abstract_elem) for abstract_elem in
abstract_elems) if abstract_elems else ""
35
36         # Store the metadata in the dictionary
37         metadata_dict[pubmed_id] = {
38             "ArticleTitle": title,
39             "AbstractText": abstract,
40             "query": query_term
41         }
42     return metadata_dict

```

Listing 7: get meta data

3.3 Overlap id

We finally retrieve PubMed article IDs for Alzheimer's and cancer, find IDs common to both, and print them. It then fetches metadata for the articles related to each disease and stores all metadata in a JSON file named metadata.json.

```

1  if __name__ == "__main__":
2      Alzheimers_id = get_id("Alzheimers")
3      Cancer_id = get_id("cancer")
4
5      common_ids = list(set(Alzheimers_id) & set(Cancer_id))
6      print("Common IDs:", common_ids)
7
8      Alzheimers_metadata = get_metadata(Alzheimers_id, "Alzheimers")
9      Cancer_metadata = get_metadata(Cancer_id, "cancer")
10
11     all_metadata = {**Alzheimers_metadata, **Cancer_metadata}
12     with open('metadata.json', 'w') as f:
13         json.dump(all_metadata, f, indent=4)

```

Listing 8: Implementation

4 Exercise 4

We install the pretrained model from huggingface.

```
1 from transformers import AutoTokenizer, AutoModel
2
3 # load model and tokenizer
4 tokenizer = AutoTokenizer.from_pretrained('allenai/specter')
5 model = AutoModel.from_pretrained('allenai/specter')
```

Listing 9: Install pretrained model

Then we extract data from the file we generated before and encode the text data using the pretrained model. We can check the dimension of the embedding vectors `embeddings[0].size` and find that the dimension is 768.

```
1 import json
2 with open('metadata.json') as f:
3     papers = json.load(f)
4 import tqdm
5
6 # we can use a persistent dictionary (via shelve) so we can stop and restart if needed
7 # alternatively, do the same but with embeddings starting as an empty dictionary
8 embeddings = {}
9 for pmid, paper in tqdm.tqdm(papers.items()):
10     data = [paper["ArticleTitle"] + tokenizer.sep_token + " ".join(paper["AbstractText"])]
11     inputs = tokenizer(
12         data, padding=True, truncation=True, return_tensors="pt", max_length=512
13     )
14     result = model(**inputs)
15     # take the first token in the batch as the embedding
16     embeddings[pmid] = result.last_hidden_state[:, 0, :].detach().numpy()[0]
17
18 # turn our dictionary into a list
19 embeddings = [embeddings[pmid] for pmid in papers.keys()]
```

Listing 10: Encode data

Then, we use PCA to reduce the high-dimension vectors to low-dimension. This helps us better interpret and visualize data. We extract three PCs and make a visualization ??.

```
1 from sklearn import decomposition
2 import pandas as pd
3
4 pca = decomposition.PCA(n_components=3)
5 embeddings_pca = pd.DataFrame(
6     pca.fit_transform(embeddings),
7     columns=['PC0', 'PC1', 'PC2']
8 )
9 embeddings_pca["query"] = [paper["query"] for paper in papers.values()]
```

Listing 11: PCA implementation

One of the ideas of doing this visualization is to see if these vectors can well-represent the articles from different categories. We find that PCA plots reveal some degree of separation between the Alzheimer's (red) and cancer (blue) especially in PC0 and PC1 or PC0 and PC2. However, there is a noticeable overlapping, especially in the PC1 vs PC2 plot, suggesting that while the SPECTER embeddings capture meaningful differences between the two categories. One explanation is that the first two PCs capture the most variance. Variance measures how spread out the data points are from the mean. In high-variance directions, the data points are widely distributed, indicating that these directions capture the underlying differences between data points more effectively. Higher variance therefore often indicates that there are more distinguishing features along that direction.

	PC0	PC1	PC2	query
0	-4.603137	-2.170447	-1.526437	Alzheimers
1	-7.595276	0.435876	-1.645755	Alzheimers
2	-6.625034	-2.008937	-1.253660	Alzheimers
3	-5.550500	2.311249	-4.340761	Alzheimers
4	-5.913295	3.007254	-0.234550	Alzheimers
...
1992	0.697405	-3.691629	5.893336	cancer
1993	3.340729	-1.148313	6.953101	cancer
1994	9.356931	-1.196874	-3.162101	cancer
1995	3.036748	-4.447534	6.247646	cancer
1996	6.032124	-1.736871	-1.251254	cancer

1997 rows x 4 columns

Figure 5: PCA results

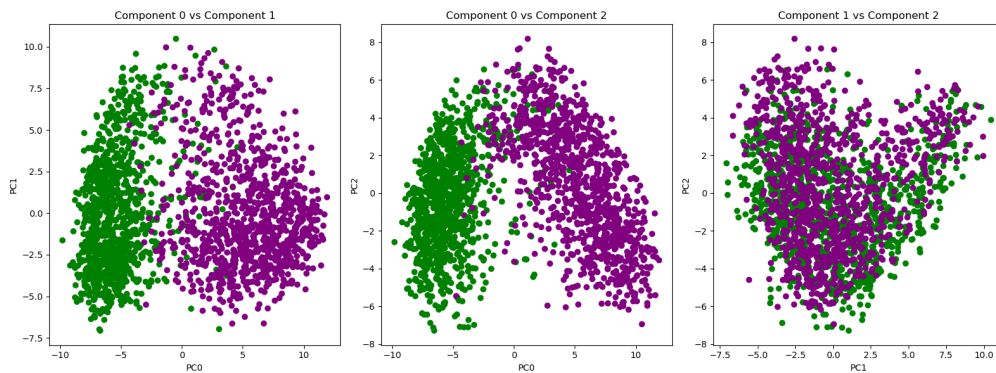


Figure 6: PCA visualization results

References