

iTerm2 Shell Edit View Session Scripts Profiles Toolbelt Window Help

IR8A.py — informaiton_retrieval

```

EXPLORER ... IR8A.py x
INFORMATON_RETRIEV... hw8 > IR8A.py > ...
> hw1
> hw2
> hw3
> hw4
> hw5
> hw6
> hw7
> hw8
> hw9
> lecture
~$HW8.docx
bible.txt
Figure_1.png
Figure_2.png
Figure_3.png
HW8.docx
IR8A.py
> hw9
> lecture
10 import re
11 import numpy as np
12 import random
13 from sklearn.feature_extraction.text import CountVectorizer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15 from sklearn.metrics.pairwise import cosine_similarity
16 from matplotlib import pyplot as plt
17 relevant = 50
18 f = open("binary.txt", "r")
19 docs = f.readlines()
20 f.close()
21
22 tfidf = TfidfVectorizer(max_df=0.4, min_df=2)
23 dt = tfidf.fit_transform(docs)
24 N = len(docs)
25 query = random.randint(0, N)
26 print(query, docs[query])
27
28 sim = cosine_similarity(dt[query], dt)
29 topNdf = set()
30 for index in np.argsort(sim)[0][-10:-10]:
31     topNdf.add(index)
32
33 print(topNdf)
34
35 cv = CountVectorizer(binary=True, max_df=0.4, min_df=2)
36 dt2 = cv.fit_transform(docs)
37 sim2 = cosine_similarity(dt2[query], dt2)
38 sorted = np.argsort(sim2)[0][-1:-1]
39 precision = np.zeros(N)
40 recall = np.zeros(N)
41 m = 0
42 for i in range(N):
43     if sorted[i] in topNdf:
44         m = m + 1
45     tp = m
46     fn = relevant - m
47     fp = i + 1 - m
48

```

ln 4, Col 35 Spaces: 4 UTF-8 CRLF Python 3.9.6 64-bit

iTerm2 Shell Edit View Session Scripts Profiles Toolbelt Window Help

IR8A.py — informaiton_retrieval

```

EXPLORER ... IR8A.py x
INFORMATON_RETRIEV... hw8 > IR8A.py > ...
> hw1
> hw2
> hw3
> hw4
> hw5
> hw6
> hw7
> hw8
> hw9
> lecture
~$HW8.docx
bible.txt
Figure_1.png
Figure_2.png
Figure_3.png
HW8.docx
IR8A.py
> hw9
> lecture
48 fp = i + 1 - m
49 tn = N - tp - fn - fp
50 precision[i] = tp / (tp + fp)
51 recall[i] = tp / (tp + fn)
52
53 # correct AP calculation
54 # first append sentinel values at the end
55 mrec = np.concatenate(([0.], recall, [1.]))
56 mpre = np.concatenate(([0.], precision, [0.]))
57
58 # compute the precision envelope
59 for i in range(mpre.size - 1, 0, -1):
60     mpre[i - 1] = np.maximum(mpre[i - 1], mpre[i])
61
62 # to calculate area under PR curve, look for points
63 # where X axis (recall) changes value
64 i = np.where(mrec[i:] != mrec[i-1])[0]
65
66 # AP= AP1 + AP2+ AP3+ AP4 + ...
67 ap = np.sum((mrec[i + 1] - mrec[i]) * mpre[i + 1])
68
69 print(f'MAP: {ap}')
70
71 plt.figure()
72 plt.title("precision-recall graph")
73 plt.scatter(recall, precision)
74
75 eleven_recalls = np.zeros(11)
76 interpolated = np.zeros(11)
77 n = 0
78 for i in range(N):
79     if n <= 10 and recall[i] * 10 >= n:
80         interpolated[n] = max(precision[i:])
81         eleven_recalls[n] = recall[i]
82         n += 1
83     if n > 10: break
84 plt.figure()
85 plt.title("eleven-point interpolated precision-recall graph")
86 plt.scatter(eleven_recalls, interpolated)

```

ln 76, Col 28 Spaces: 4 UTF-8 CRLF Python 3.9.6 64-bit

iTerm2 Shell Edit View Session Scripts Profiles Toolbar Window Help

IR8A.py — informaiton_retrieval

```

EXPLORER ... IR8A.py
hw1
hw2
hw3
hw4
hw5
hw6
hw7
hw8
~$HW8.docx
bible.txt
Figure_1.png
Figure_2.png
Figure_3.png
HW8.docx
IR8A.py
hw9
lecture

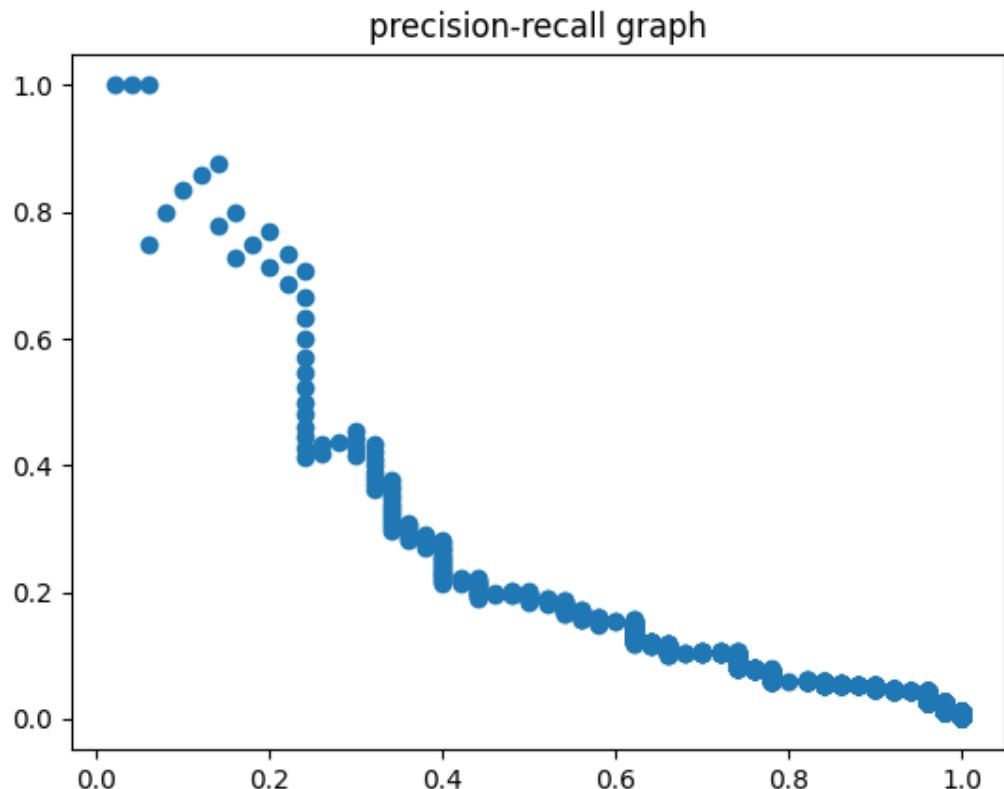
```

```

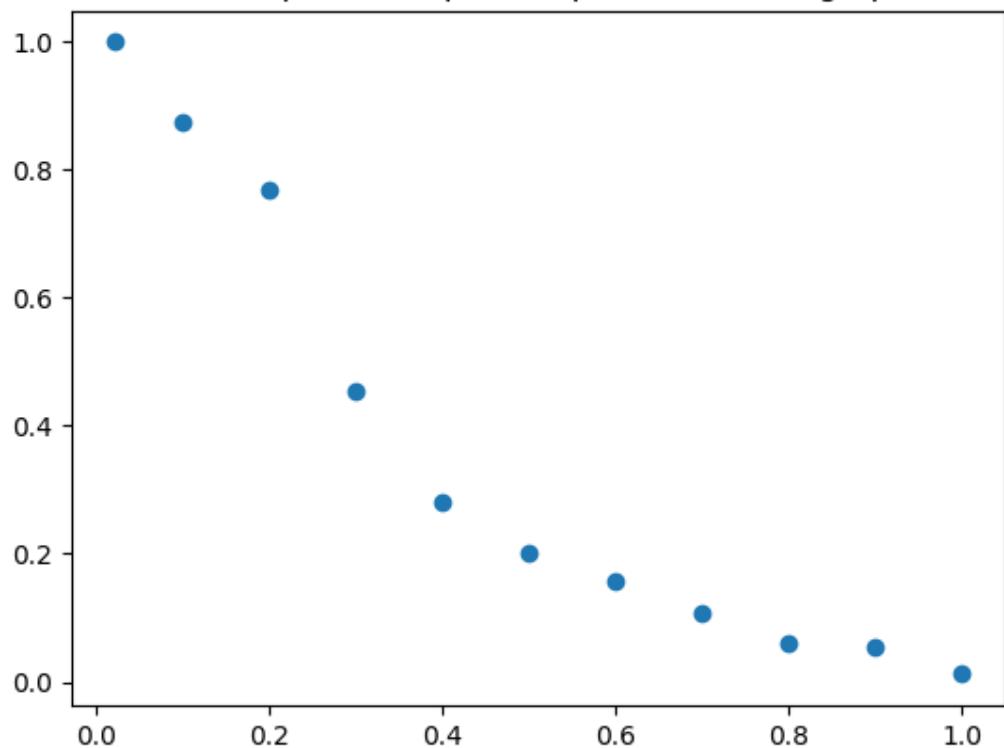
67 ap = np.sum((mrec[i + 1] - mrec[i]) * mpre[i + 1])
68
69 print(f'MAP: {ap}')
70
71 plt.figure()
72 plt.title("precision-recall graph")
73 plt.scatter(recall, precision)
74
75 eleven_recalls = np.zeros(11)
76 interpolated = np.zeros(11)
77 n = 0
78 for i in range(N):
79     if n <= 10 and recall[i] * 10 >= n:
80         interpolated[n] = max(precision[i:])
81         eleven_recalls[n] = recall[i]
82         n += 1
83     if n > 10: break
84 plt.figure()
85 plt.title("eleven-point interpolated precision-recall graph")
86 plt.scatter(eleven_recalls, interpolated)
87
88 rocx = np.zeros(N)
89 recall = np.zeros(N)
90 m = 0
91 for i in range(N):
92     if sorted[i] in toptfidf:
93         m += 1
94         tp = m
95         fn = relevant - m
96         fp = i + 1 - m
97         tn = N - tp - fn - fp
98         rocx[i] = tp / (tp + fn)
99         recall[i] = tp / (tp + fn)
100 plt.figure()
101 plt.title("ROC curve")
102 plt.scatter(rocx, recall)
103
104 plt.show()

```

Last login: Tue Sep 20 12:57:41 on ttys000
(base) jbchang@jbchangs-mbp:~/informaiton_retrieval/hw8
(base) jbchang@jbchangs-mbp:~/informaiton_retrieval/hw8 python3 IR8A.py
8507 And the King said unto him, Turn aside, and stand here. And he turned aside
, and stood still.
{29058, 3972, 8669, 8071, 5383, 7190, 27833, 21666, 15395, 30371, 9253, 95
11, 17833, 24489, 1581, 1582, 28982, 7479, 7480, 25914, 8507, 19645, 23104, 6465
, 22471, 7371, 14796, 26317, 4430, 21715, 19669, 17494, 6616, 3803, 5084, 23776,
17598, 9446, 5223, 3829, 17132, 29680, 13889, 14450, 5235, 12276, 7035, 7038, 1
6895}
MAP: 0.11585215760555412
(base) jbchang@jbchangs-mbp:~/informaiton_retrieval/hw8



eleven-point interpolated precision-recall graph



ROC curve

