

به نام خدا



گزارش پروژه منطق فازی

درس هوش محاسباتی

دانشگاه اصفهان

دانشکده مهندسی کامپیوتر

استاد درس : دکتر حسین کارشناس

تهیه کنندگان :

علیرضا دستمالچی ساعی، محمد حسین دهقانی

خرداد ۱۴۰۲

مقدمه

در این پروژه خواسته شده بود که از یک دیتاست مخصوص پیام‌های ایمیل که مشخص می‌کند ایمیل اسپم است یا نه، یک سیستم استنتاج فازی با یک مجموعه قانون بسازیم که با گرفتن ویژگی‌های استخراج شده از متن یک ایمیل، اسپم بودن یک ایمیل را تشخیص دهیم.

پیاده‌سازی

این برنامه به دنبال پیدا کردن بهترین مدل رده‌بندی اسپم و غیر اسپم برای داده‌های متنی است. برای این کار از الگوریتم ژنتیک استفاده شده است.

در این برنامه، ابتدا داده‌های متنی از یک فایل خوانده شده و سپس آن‌ها به یک دیتافریم تبدیل می‌شوند. سپس از روش TF-IDF برای استخراج ویژگی‌های مهم استفاده می‌شود و سپس با استفاده از روش کاهش بعد PCA، تعداد ویژگی‌ها (تعداد بعدها) کاهش داده می‌شود. در ادامه، با استفاده از روش feature_selection تنها ویژگی‌های مهم برای دسته‌بندی انتخاب می‌شوند.

الگوریتم تکاملی

برای الگوریتم تکاملی، از کتابخانه آماده Deap که یک کتابخانه قوی و جامع برای اجرای الگوریتم‌های تکاملی است، استفاده کرده‌ایم. ابتدا ژنوتایپ یا کروموزوم (Individual) و تابع ارزیابی (Evaluation) را برای toolbox کتابخانه تعریف می‌کنیم.

نمایش ما (ژنوتایپ) به صورت یک دیکشنری می باشد که در toolbox رجیستر شده و دارای ۲ بخش می باشد:

- قسمت Rules که شامل لیستی از قوانین است که به آن Rule base نیز می گوین، بدین صورت که هر قانون X_1 تا X_n و Y دارد. تعداد قوانین هم بین ۵۰ تا ۵۰۰ عدد می باشد.

- قسمت Linguistic Functions که برای هر متغیر زبانی یک تابع درجه عضویت تعریف کرده ایم که در ابتدا به صورت رندوم انتخاب می شود.

سپس، برای تغییر، از توابعی که خودمان تعریف کرده ایم استفاده می کنیم:

- **بازترکیب:**

برای بازترکیب از عملگر تغییر، از one_point_crossover استفاده کرده ایم.

- **جهش:**

برای جهش از عملگر تغییر، با یک احتمالی مقادیر قانون ها را تغییر می دهیم سپس با همان احتمال قبلی، توابع عضویت مخصوص هر مقدار زبانی، m و s را جهش می دهیم.

برای تابع Evaluate که برازندگی هر ژنوتایپ (کروموزم) را حساب می‌کند، با استفاده فرمول‌هایی که در داک ذکر شده‌اند، استفاده شده است.

نتایج

برنامه برای تعداد جمعیت ۱۰ و تعداد نسل‌های ۱۰ با ۳ مقدار زبانی، نتیجه زیر را با دقت 90 حاصل کرده است (۱۱ قانون):

Gen	nevals	avg	std	min	max
۰	10	2728.68	2523.3	0	7762.16
۱	8	5190.56	1879.14	3164.81	9882.28
۲	9	6383.23	1980.43	4108.52	9882.28
۳	9	6314.78	2523.94	2983.26	9882.28
۴	9	5409.89	3965.59	0	10636.3
۵	9	6222.58	3907.63	596.664	12346.7
۶	9	7305.13	3014.4	328	9761.59
۷	10	5834.59	4636.2	0	14877
۸	9	6857.43	4724.04	0	16916
۹	9	12228.9	4046.9	4330.8	16923.6
۱۰	8	12447.1	4932.68	4096.39	16923.6

Best individual: {'rules': [[1, 0, 2, 1, 0, 0], [2, 0, 1, 2, 1, 1], [2, 2, 2, 1, 0, 1], [2, 0, 0, 0, 0, 0], [0, 0, 1, 2, 1, 1], [1, 0, 1, 2, 1, 1], [0, 2, 2, 2, 1, 1], [0, 1, 2, 2, 1, 1], [2, 0, 1, 2, 1, 1], [0, 0, 0, 2, 0, 1], [2, 0, 2, 0, 1, 0]], 'ling_funcs': {0: (<function rect_trap at 0x7f7b6cd1b400>, 0.2597256470909224, 0.24841580892219772), 1: (<function rect_trap at 0x7f7b6cd1b400>, -0.053399118156442604, 0.9552181094331375), 2: (<function gaussian at 0x7f7b6cd1ae60>, -0.4030993382928567, 0.4065217311866878)}}}

accuracy 0.9063509149623251

- با تحلیل نتایج می‌توان به موارد زیر دست یافت:
- علت برخی min هایی که در وسط اجرا صفر شده‌اند به دلیل mutation می‌باشد.
 - الگوریتم تکاملی روند صعودی را طی کرده است بجز در نسل ۶ که افت کوچکی داشته اما سپس پیشرفت نسبتاً خوبی داشته است.
 - بهترین ژنوتایپ ما که در بالا آورده شده است، برای ۳ مقدار زبانی با توابع درجه عضویت زیر است:
 - Low: rect_trap ○
 - Medium: rect rap ○
 - High: gaussian ○

با تحلیل قوانین بدست آمده در بهترین ژنوتایپ با اطلاعات ذکر شده در بالا می‌توان فهمید که اکثر ایمیل‌هایی که اسپم تشخیص داده شدند مقادیر زبانی

Medium یا High را در ویژگی‌ها داشتند به خصوص ویژگی چهارم که به تنهایی در یک قانون باعث اسپم تشخیص داده شدن آن ایمیل شده است.

در نهایت بدست آوردن همچنین نتیجه‌ای با تقسیم ۵۰ درصدی آموزش و تست عملکرد خوب الگوریتم تکاملی و ژنوتایپ بدست آمده دارد.

در مرحله‌ی دوم با اجرای این کد با تعداد نسل‌های ۱۰ و تعداد جمعیت ۱۰ ولی با ۵ مقدار زبانی و ۳۰ قانون نتایج زیر حاصل شده‌است:

Gen	nevals	avg	std	min	max
۰	10	8831.64	7137.22	0	21031.6
۱	7	13882.7	6875.92	891.316	21882.6
۲	9	13874.7	10396.3	685.361	34121.3
۳	7	13236.3	11459.6	1787.57	34121.3
۴	9	15715.8	8394.14	4913.41	34121.3
۵	10	15381.1	6592.33	7540.41	29822.3
۶	10	18054.3	8689.16	1902.88	28923.3
۷	10	24218.1	4993.08	14565.1	29614.1
۸	9	20289.8	9869.81	740.027	29614.1
۹	10	18864.4	10102.3	0	29614.1
۱۰	10	20078.8	9755.71	1958.4	32511

Best individual: {'rules': [[4, 3, 1, 4, 1, 1], [2, 4, 4, 2, 3, 0], [1, 1, 1, 3, 4, 1], [0, 1, 3, 3, 2, 1], [4, 2, 0, 1, 3, 1], [2, 0, 3, 3, 1, 1], [1, 1, 2, 4, 4, 0], [3, 3, 0, 2, 1, 1], [3, 2, 3, 2, 0, 1], [1, 4, 1, 1, 0, 1], [2, 1, 1, 4, 1, 0], [2, 2, 1, 2, 3, 0], [3, 2, 2, 0, 1, 1], [3, 1, 3, 4, 1, 1], [4, 2, 2, 0, 3, 0], [3, 0, 0, 2, 0, 1], [1, 1, 2, 4, 1, 1], [3, 4, 4, 1, 2, 0], [0, 2, 3, 4, 0, 0], [1, 3, 2, 2, 4, 0], [2, 0, 2, 1, 3, 1], [0, 4, 4, 0, 4, 0], [0, 3,

```
4, 1, 2, 1], [3, 0, 4, 4, 1, 1], [4, 4, 1, 4, 1, 1], [1, 4, 1, 3, 2, 1], [1, 4,
0, 4, 4, 0], [4, 3, 0, 4, 1, 0], [4, 1, 1, 4, 1, 0], [4, 3, 4, 4, 1, 1]],
'ling_funcs': {0: (<function sigmoid at 0x7f7b6cd1ad40>,
0.16739159164262007, 0.325079136985742), 1: (<function
rect_trap at 0x7f7b6cd1b400>, -0.17868170411156825,
0.6728158429298997), 2: (<function rect_trap at
0x7f7b6cd1b400>, -0.7163853770078319,
0.29962110739064896), 3: (<function rect_trap at
0x7f7b6cd1b400>, 0.8200094237933757,
0.7496362214724762), 4: (<function sigmoid at
0x7f7b6cd1ad40>, 0.9151647467099906,
0.8687030234852108))}}
accuracy 0.7192825112107624
```

با تحلیل نتایج بدست آمده، می‌توان به موارد زیر رسید:

- ژنوتایپ‌ها در الگوریتم به مقداری کاهش برازندگی داشته‌اند به طوریکه ماکسیمم ۳۴ هزار بوده ولی بهترین ژنوتایپ برازندگی ۳۲ هزار دارد.
- صفرهای موجود در وسط الگوریتم تکاملی بدلیل mutation می‌باشد.
- توابع عضویت برای مقادیر زبانی به صورت زیر می‌باشد:

- Very low: sigmoid ○
- Low: rect_trap ○
- Medium: rect_trap ○
- High: rect_teap ○
- Very high: Sigmoid ○

با افزایش تعداد قوانین و مقادیر زبانی با ۱۰ نسل و تست ۴۰ درصد به دقت ۷۱ رسیدیم که نسبت به حالت قبلی مقداری افت دقت داشته‌ایم و ممکن به است دلایل مختلفی داشته باشد، مانند:

- مقداردهی اولیه به صورت رندوم با مقادیر خاصی شروع کرده است
- میزان Exploration با احتمالات داده شده برای این الگوریتم کم بوده
- به دلیل علت قبلی، در نقطه بهینه محلی گیر کرده‌ایم.

در مرحله سوم نیز، با اجرای کد به تعداد ۱۰ نسل و تعداد جمعیت ۱۰ ولی با تعداد قوانین ۵۰ و تعداد متغیرهای زبانی ۷، نتایج زیر حاصل شده است:

Gen	nevals	avg	std	min	max
۰	10	8647.64	4569.4	3369.8	17920.3
۱	10	18172.4	7971.29	4556.83	29442.2
۲	10	24970.3	9514.38	6760.4	44509.3
۳	6	26287	5022.06	18944.7	36058.8
۴	10	29019	11429.2	12630.7	41783.6
۵	10	26321.3	17298.2	789.168	53192.8
۶	10	28479.5	14005.3	6863.96	50048.8
۷	8	32525.9	17258.7	5062.54	52320.3
۸	10	45356.9	7693.34	25726.7	52320.3
۹	9	42068.1	16809.6	6172.81	55757.2
۱۰	8	34474.9	20993.1	725.634	55795

Best individual: {'rules': [[3, 5, 2, 2, 4, 1], [5, 6, 1, 3, 1, 0], [1, 1, 5, 6, 4, 0], [2, 4, 3, 3, 3, 0], [3, 5, 1, 6, 1, 0], [2, 3, 4, 2, 2, 0], [5, 3, 4, 4, 0, 0], [1, 3, 0, 3, 3, 1], [3, 3, 0, 2, 4, 1], [3, 1, 4, 0, 1, 1], [2, 3, 2, 6, 6, 0], [2, 3, 4, 4, 4, 0], [3, 4, 1, 6, 1, 0], [3, 3, 5, 4, 5, 1], [5, 5, 2, 1, 0, 1], [2, 2, 5, 5, 2, 1], [4, 3, 0, 3, 5, 1], [1, 5, 2, 0, 5, 0], [1, 1, 2, 4, 5, 1], [2, 3, 0, 0, 0, 1], [5, 4, 6, 6, 1, 1], [1, 2, 0, 6, 1, 0], [5, 3,

3, 0, 0, 0], [4, 0, 3, 1, 2, 1], [1, 4, 0, 1, 3, 1], [3, 4, 1, 0, 0, 1], [3, 5, 2, 1, 4, 1], [3, 4, 6, 5, 4, 0], [4, 3, 3, 1, 2, 0], [4, 1, 1, 4, 5, 1], [5, 1, 2, 2, 4, 1], [4, 5, 3, 6, 5, 1], [4, 1, 2, 2, 6, 0], [5, 6, 1, 4, 0, 1], [2, 5, 0, 2, 5, 1], [5, 1, 0, 1, 1, 0], [4, 5, 6, 4, 0, 0], [3, 5, 2, 0, 6, 0], [1, 2, 1, 0, 1, 0], [4, 2, 0, 1, 5, 0], [1, 2, 0, 4, 5, 0], [0, 2, 3, 6, 4, 0], [5, 0, 2, 1, 5, 1], [4, 6, 5, 1, 4, 1], [4, 0, 6, 5, 3, 1], [2, 0, 3, 0, 4, 1], [0, 0, 1, 1, 2, 0], [0, 3, 6, 2, 2, 0], [2, 5, 5, 6, 0, 1], [4, 6, 1, 2, 3, 0]],

'ling_funcs': {0: (<function sigmoid at 0x7f7b6cd1ad40>, 0.4882531421828613, 0.9252488574391238), 1: (<function sigmoid at 0x7f7b6cd1ad40>, 0.18871791792197468, 0.08519210679339584), 2: (<function gaussian at 0x7f7b6cd1ae60>, 0.053283191970501464, 0.7669474526119026), 3: (<function sigmoid at 0x7f7b6cd1ad40>, 0.867398853556133, 0.6941425929969136), 4: (<function rect_trap at 0x7f7b6cd1b400>, 0.8327768822356998, 0.6584981422332299), 5: (<function sigmoid at 0x7f7b6cd1ad40>, 0.9729642115566914, 0.12098425766441023), 6: (<function sigmoid at 0x7f7b6cd1ad40>, 0.6309064453410838, 0.8866333122999813))}}

accuracy 0.7726457399103139

با تحلیل نتایج بدست آمده، می‌توان به موارد زیر رسید:

- ژنوتایپ‌ها در الگوریتم پسرفت نداشته به صورت افزایشی برازندگی آنها آپدیت می‌شد.

- توابع عضویت برای مقادیر زبانی به صورت زیر می‌باشد:

 - Very very low: sigmoid

 - Very low: sigmoid

 - Low: gaussian

 - Medium: sigmoid

 - High: rect_trap

 - Very high: sigmoid

 - High: Sigmoid

با افزایش تعداد قوانین و مقادیر زبانی با ۱۰ نسل و تست ۴۰ درصد به دقت ۸۷ رسیدیم که نسبت به حالت قبلی افزایش دقت ۶ درصدی حاصل شد.

یکی از نتایج دیگر با ۲۰ قانون و تنظیمات اجرای قبلی:

```
CI_HW3.ipynb ☆
Fichier Modifier Affichage Insérer Exécution Outils Aide

+ Code + Texte

print("Best individual: ", best_individual)
prediction = []
for x in X_test:
    prediction.append(y_hat(x, best_individual['rules'], best_individual['ling_funcs']))
corrects = 0
for i in range(len(y_test)):
    if y_test[i] == prediction[i]:
        corrects += 1
print('accuracy', corrects/len(y_test))

gen    nevals    avg      std      min      max
0      10         7643.05  4562.87  0         15560.7
1      8         7381.19  2406.34  3807.45   11018.5
2      9         12026.3  4733.95  1787.19   19677.5
3      8         11486.4  2849.03  6026      15519.5
4      10        11573.9  3937.55  4468.77   15939
5      10        10638.6  4754.33  1171.14   18088.1
6      10        11059.4  4250.18  893.99    14789
7      8         13018.2  4533.11  5693.25   19594.7
8      7         16330.9  3037.1   8636.84   19594.7
9      7         16281.3  3727.34  8035.91   19925.8
10     10        15947.6  6529.22  440.134   20017.9

Best individual: {'rules': [[1, 0, 1, 1, 2, 0], [2, 1, 1, 2, 1, 1], [2, 2, 1, 0, 1, 0], [0, 1, 1, 0, 0, 0], [1, 2, 1, 1, 2, 0], [0, 1, 2, 2, 2, 0], [2, 1, 0, 2, 1, 0], [1, 1, 1, 0, 0, 0],
accuracy 0.910762331838565
```

یکی از نتایج دیگر با ۳۰ قانون و تنظیمات اجرای قبلی:

```

CJ_HW3.ipynb
Fichier Modifier Affichage Insérer Exécution Outils Aide

+ Code + Texte

print("Best individual: ", best_individual)
prediction = []
for x in X_test:
    prediction.append(y_hat(x, best_individual['rules'], best_individual['ling_funcs']))
corrects = 0
for i in range(len(y_test)):
    if y_test[i] == prediction[i]:
        corrects += 1
print('accuracy', corrects/len(y_test))

gen  nevals  avg      std      min      max
0    10      7675.09  6434.28  -0       19147.1
1    9       11471.1  5848.92  441.807  22284.3
2    9       15255.1  6238.04  446.014  22921.8
3    9       15705.9  7770.61  447      27928
4    7       23478.4  4311.39  15888.5  32314.7
5    10      17534.6  9803.71  0        31706
6    9       20801.2  11035.6  -0       31706
7    10      21841   9332.51  1451.8   32170.3
8    10      28581.9  4003.93  21026.7  32170.3
9    10      25168.9  10154.2  -0       34683.1
10   8       31195   3825.19  20615.4  34683.1

Best individual: ('rules': [[1, 2, 0, 0, 2, 1], [2, 0, 0, 1, 2, 0], [2, 1, 1, 2, 1, 1], [1, 2, 2, 1, 0, 1], [2, 0, 0, 0, 2, 1], [2, 0, 2, 1, 0, 0], [0, 0, 2, 2, 1, 1], [2, 0, 2, 1, 2, 2],
accuracy 0.7896860986547085

```

استفاده از روش‌های کاهش بعد (PCA) و انتخاب ویژگی‌ها (feature selection) بستگی به مجموعه داده و مشکل خاصی دارد که به آن پرداخته می‌شود. استفاده از این تکنیک‌ها در این پروژه باعث شد که عملکرد مدل ما با کاهش تعداد ویژگی‌ها و تمرکز بر مرتبط‌ترین آنها بهبود بخشد و از محاسبات اضافی راحت شدیم. برای بهتر شدن نتیجه‌ها شاید می‌شد که از روش‌های دیگر کاهش بعد مانند MDA و ... نیز استفاده کنیم تا تاثیر دیگر روش‌ها را دیده و بهترین نتیجه را انتخاب کنیم.