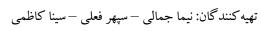
بازيابي پيشرفته اطلاعات

مدرس: دکتر بیگی شماره گروه: ۵





گزارش فاز دوم پروژه

فهرست مطالب

ساختن فضای برداری برای دستهبندها	٣
پیش پردازش و تنظیم فیلدهای اولیه	٣
پیادهسازی فضای برداری tf-idf به روش ntn	۴
ذخیرهسازی فضای برداری	۵
پیادهسازی دستهبندها	۶
Naïve Bayes	9
k-NN	٨
SVM	4
Random Forest	1.
محاسبهی Validation Set	11
محاسبهی بهترین پارامتر برای الگوریتم k-NN	14
محاسبهی بهترین پارامتر برای الگوریتم SVM	١٣
بهبود سیستم بازیابی فاز اول پروژه	14
تابع اجماع نظر دستهبندها برای تعیین برچسب هر مستند	10
جستجو بر اساس دسته	١۵
ارزیابی نهایی	17
پیادهسازی تابع find_metric	14
معیارهای ارزیابی برای Naïve Bayes	14
معیارهای ارزیابی برای k-NN	1.4
معیارهای ارزیابی برای SVM	1.4
معیارهای ارزیابی برای Random Forest	1.4
نتیجه گیری نهایی معیارهای ارزیابی	19

Y•	حوهى تقسيم وظايف
Y•	يما جمالى
Υ.	سپهر فعلی
۲.	سينا كاظمى
Y1	براجع

ساختن فضای برداری برای دستهبندها

در این بخش که در واقع بخش آغازین کار است، ابتدا کلاس جدیدی به نام Classifier تعریف می کنیم و IRSystem تعریف می کنیم. همچنین به کلاس IRSystem که در فاز اول پروژه طراحی شده بود، متدی برای محاسبه ی فضای برداری اسناد به روش tf-idf اضافه می کنیم. در ادامه به شرح پیاده سازی این بخش می پردازیم.

پیش پردازش و تنظیم فیلدهای اولیه

کلاس path یک classifier یک path را به عنوان ورودی دریافت می کند. این path مشخص می کند که دسته بندی قرار است بر روی چه فایلی اعمال شود. سپس از کلاس IRSystem برای پیش پردازش مستندات موجود در train.csv بر روی چه فایلی اعمال شود. سپس از کلاس call_create_positional برای این مجموعه مستندات صدا می زنیم تا نمایه ی می کنیم. تابع positional و call_prepare براهیم مستندات موجود در path را هم در این مجموعه درج کنیم، سایز مجموعهی آموزش را در متغیری به نام train_size نگه می داریم. سپس مستندات موجود در path را درج می کنیم. مجموعهی آموزش را در متغیری به نام train_size نگه می داریم. سپس مستندات موجود در path و path همان برای محاسبه ی معیارهای ارزیابی استفاده می شود، تنها در صورتی تعریف می کنیم که path برابر با لیستی خواهد بود که مقادیر آن، به ترتیب برابر با مقادیر ستون views در داده ی آزمون است.

```
class Classifier:
    def __init__(self, path):
        self.train_ir_sys = IRSystem(["description", "title"], "data/train.csv", None)
        self.train_ir_sys.call_prepare("english", False)
        self.train_ir_sys.call_create_positional("english")
        self.train_ir_sys.csv_lnsert(path, "english")
        self.train_ir_sys.csv_lnsert(path, "english")
        self.train_ir_sys.csv_lnsert(path, "english")
        self.train_ir_sys.csv_lnsert(path, "english")
        self.y_train = self.csv_views("data/train.csv")

        self.y_train = self.csv_views("data/train.csv")

        self.vector_space = self.train_ir_sys.use_ntn("english")

        self.y_test = None
        self.vave_classifier = None
        self.naive_bayes_classifier = None
        self.naive_bayes_classifier = None
        if path == "data/test.csv":
            self.y_test = self.csv_views("data/test.csv")

def csv_views(self, path):
        df = pd.read_csv(path, usecols=["views"])
        result = []
        for i in range(len(df)):
            result += [df.iloc(i]["views"]]
        return result
```

تصویر ۱: طراحی اولیه کلاس Classifier و تنظیم فیلدهای آن

پیادهسازی فضای برداری tf-idf به روش ntn

در این بخش، ابتدا در کلاس IRSystem متدی با نام ntn طراحی کردیم که یک doc_id و زبان آن را ورودی می گیرد. سپس یک دیکشنری به نام result می سازد که کلیدهای آن، ترمهای آن مستند و مقادیر آن، حاصل عبارت می گیرد. سپس یک دیکشنری به نام tresult می این تابع توسط تابعی به نام use_ntn صدا زده می شود. این تابع به ازای تمام مستندات، دیکشنری گفته شده را محاسبه کرده و بر می گرداند.

```
def ntn(self, doc_id, lang):
    doc = []
    for part in self.structured_documents[lang][doc_id]:
        for word in part:
            doc += [word]
    doc_dict = Counter(doc)
    result = dict()
    for term in doc_dict.keys():
        if term not in result.keys():
            p = self.positional_index[lang][term]
            df = len(p.kevs()) - 1
            idf = math.log10(len(self.structured_documents[lang]) / df)
            result[term] = doc_dict[term] * idf
    return result
def use_ntn(self, lang):
   vectors = []
   for doc_id in range(len(self.structured_documents[lang])):
       vectors += [self.ntn(doc_id, lang)]
```

تصویر ۲: طراحی توابع ntn و use_ntn

برای طراحی فضای برداری که الگوریتمهای کتابخانههای آماده برای SVM و Random Forest بتوانند آن را پیاده کنند، نیاز به ماتریسی از pndarrayها داریم و نمی توان از دیکشنری استفاده کرد. به همین دلیل نیاز داریم برای هر ترم مستقل در کل دیکشنری (train + test) یک عدد داشته باشیم که به صورت یکتا بین آنها رابطه وجود داشته باشد. برای این کار تابع token_to_number را طراحی می کنیم که به هر token عدد خاصی نسبت می دهد. این عدد خاص از صفر تا 1 – (len(terms) می تواند باشد.

سپس برای ساختن ماتریسی از ereate_vector_space را صدا می زنیم که ابتدا ماتریسی از صفر تولید می کند که تعداد سطرهای آن برابر با تعداد مستندات و تعداد ستونهای آن برابر با تعداد کل ترمهاست. در نهایت با صدا زدن use_ntn برای هر مستند، برای هر ترم در آن مستند از تابع token_to_number استفاده می کنیم

و وزن مربوط به آن ترم در مستند را در ماتریس بهروز رسانی می کنیم. در نهایت این ماتریس در فیلدی به نام train_vector_space ذخیره می شود.

```
def token_to_number(self, ir_sys, lang):
    return {token: ind for ind, token in enumerate(ir_sys.positional_index[lang].keys())}

def create_vector_matrix(self, ir_sys, lang):
    vector = np.zeros([len(ir_sys.structured_documents[lang]), len(ir_sys.positional_index[lang].keys())])
    tokenized_vector = ir_sys.use_ntn(lang)
    for doc_id in range(len(tokenized_vector)):
        for term in tokenized_vector[doc_id].keys():
            vector[doc_id][self.token_to_number(ir_sys, lang)[term]] = tokenized_vector[doc_id][term]
    return vector
```

تصویر ۳: طراحی توابع مورد نیاز برای ساخت فضای برداری به روش ماتریسی

ذخیرهسازی فضای برداری

فضای برداری ایجاد شده در قسمت قبل بسیار بزرگ است، زیرا تعداد ترمهای دیکشنری چیزی در حدود ۱۲۰۰۰ است! لذا هر نمونه از Classifier که ساخته می شود، زمان زیادی می گیرد. برای همین تصمیم گرفتیم دو نمونهای که از کلاس Classifier می سازیم (یک نمونه برای ted_talks.csv و یک نمونه برای ted_talks.csv) را ذخیره کنیم، البته این ذخیره سازی بعد از تکمیل سایر متدهای کلاس Classifier انجام می گیرد. با دستور create می توانیم یک نمونه ایجاد کنیم و با دستور save آن نمونه را ذخیره کنیم. نمونه ی این دستورات را می توان در این تست مشاهده کرد.

پیادهسازی دستهبندها

در این بخش توابع دستهبندی را طراحی می کنیم. همچنین توابع مورد نیاز برای محاسبه ی بهترین پارامتر را طراحی می کنیم. در ادامه به شرح پیاده سازی این توابع می پردازیم. دقت شود که از کل مجموعه داده ی train.csv برای آموزش دستهبندها استفاده شده است.

Naïve Bayes

یکی از راه های پیش بینی کردن گروه و دسته یک شی (در این پروژه، همان مستند)، الگوریتم Naive-Bayes است. این الگوریتم، تعداد نسبتا زیادی داده تحت عنوان داده ی آموزش دریافت می کند که گروه هریک از این داده ها مشخص است و با پردازش آن، به نوعی ویژگی های هر گروه معلوم می شود. بدین شکل که اگر x رخ داده باشد، احتمال این که به گروه A مرتبط باشد بیشتر است یا گروه B. برای پیاده سازی این الگوریتم از سه تابع بهره بردیم:

- تابع (naive_bayes_train(self, flag_counter, words, lang): این تابع، قسمت آموزش الگوریتم را انجام میدهد و درواقع پایه و اساس کار کرد الگوریتم را شامل می شود. اما نحوه کار:
 - متغیر flag_counter یک دیکشنری ۴عضوی است، متشکل از مقادیر زیر:
 - تعداد مستنداتی که عضو گروه ۱ هستند.(positive_docs)
 - تعداد مستنداتی که عضو گروه ۱- هستند.(negative_docs)
 - تعداد ترم های غیریکتا در مستندات گروه ۱ (positive_terms)
 - تعداد ترم های غیریکتا در مستندات گروه ۱- (negative_terms)

متغیر words یک دیکشنری بزرگ از ترمهای موجود در مجموعهی مستندات دادههای آموزش است. هر عضو آن، خود یک دیکشنری ۲عضوی است که از مقادیر زیر تشکیل شدهاست.

- تعداد تکرار کلمه موردنظر در مستندات گروه ۱ (["positive"]](words
- تعداد تكرار كلمه موردنظر در مستندات گروه ۱- (["negative"]](words[term]

این دو متغیر، در این تابع و با پردازش دادههای آموزش که قبلا پیش پردازش شدهاند و در متغیری از کلاس ir_system ذخیره شدهاند، مقداردهی و تکمیل می شوند تا در ادامه مورد استفاده قرار گیرند.

• تابع (naive_bayes_test(self, flag_counter, words, lang: این تابع، قسمت آزمایش الگوریتم را انجام می دهد و درواقع، خروجی الگوریتم است. یعنی با استفاده از نتایج بر آمده از پردازش داده های آموزش در تابع naive_bayes_train، گروه های داده های آزمایش را پیش بینی می کند. اما نحوه ی کار چگونه است:

به ازای هر مستند ورودی، به روش زیر می توانیم گروه مربوط به آن را حدس بزنیم.

$$C_{MAP} = \underset{c \in C}{\operatorname{argmax}} \left[Log P(c) + \sum_{1 \le k \le n_d} Log P(t_k|c) \right]$$

$$P(t|c) = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}, \qquad P(c) = \frac{len(C = c)}{len(C = c) + Len(C! = c)}$$

پیاده سازی این روابط در کد به این صورت است که

- برای محاسبه T_{ct} به words[t][c] برای محاسبه است.
- . برای محاسبه $\sum_{t' \in V} T_{ct'}$ به flag_counter[c_terms] رجوع می کنیم.
 - مقدار B برابر با تعداد ترمهای یکتای دادههای آزمون است.
- مراجعه می کنیم. len(C = c) مراجعه می کنیم.

در نهایت، نتیجه حدسهای خود را در لیستی به سایز دادههای آزمایش به نام y_predicted میریزیم که این متغیر، خروجی تابع خواهد بود.

• تابع (naive_bayes(self, lang: این تابع فرآیند اجرای الگوریتم را کنترل می کند. یعنی دیکشنریهای naive_bayes_train: این تابع فرآیند اجرای الگوریتم را words و سپس praive_bayes_train و سپس ابتدا تابع آموزش یعنی y_predicted را برمی گرداند.

تصویر ۴: توابع مورد نیاز برای دستهبند Naïve Bayes

در کنسول می توان با دستور naive_bayes test، بر چسبهایی که این دسته بند برای مجموعه ی تست می زند مشاهده کرد.

k-NN

دسته بند k-NN بسته به پارامتر k ورودی، k نزدیک ترین همسایه ی مستند داده شده را در مجموعه ی آموزش می یابد و با توجه به اکثریت برچسب ها در میان آن k مستند، برچسب مستند آزمایشی را محاسبه می کند.

به این منظور تابعی به نام knn طراحی می کنیم که مجموعه ی مستندات آموزش و برچسبهای آنها، مجموعه ی مستندات آزمایش و پارامتر k را ورودی می گیرد. می توانستیم از فضای برداری که در بخش پیاده سازی فضای برداری t در بخش توانستیم از فضای برداری t در بخش پیاده سازی فضای برداری t در بخش با آن روش به علت این که صفرهای زیادی را بدون آن که در محاسبه نیاز شوند، نگه می داریم و در واقع ما تریس sparse است، زمان زیادی می گیرد.

به همین علت فضای برداری جدیدی با نام vector_space (یا به عبارت بهتر knn_vector_space) ایجاد می کنیم که برای محاسبه ی آن، تابع use_ntn را صدا می زنیم. سپس باید فاصله ی دو مستند را حساب کنیم. برای این کار تابعی به نام two_doc_distance تعریف می کنیم که دو مستند را از فضای برداری vector_space می گیرد و فاصله ی آنها را محاسبه می کند. برای محاسبه ی لم نزدیک ترین همسایه، باید فاصله ی هر مستند آزمایش با هر مستند آموزشی را داشته باشیم که این کار از طریق تابع documents_distances انجام می شود و خروجی آن ماتریسی است که سطرهای آن، مستندات آموزش و ستونهای آن، مستندات آزمایش است.

تصویر ۵: توابع محاسبه کنندهی فاصلهی مستندات

در تابع knn، ابتدا این ماتریس محاسبه می شود. سپس هر ستون را که متناظر یک مجموعه ی تست است، به صورت صعودی مرتب کرده و k تای اول را برای هر ستون انتخاب می کنیم. سپس برچسب هر یک از این k مستند را بررسی می کنیم. اگر تعداد ۱ ها بیشتر باشد، برچسب مستند آزمایش هم ۱ و در غیر این صورت ۱- خواهد بود.

k-NN تصویر $% = \frac{1}{2} \left(\frac{1}{2} \right)$ تصویر

در کنسول هم می توان با دستور knn test k برچسبهایی که دسته بند k-NN به داده های تست می زند، مشاهده کرد. دقت شود که k یک عدد طبیعی است که کاربر ورودی می دهد.

SVM

بردار های پشتیبان به زبان ساده، مجموعهای از نقاط در فضای n بعدی داده ها هستند که مرز دسته ها را مشخص می کنند و دسته بندی داده ها بر اساس آنها انجام می شود و با جابه جایی یکی از آنها خروجی دسته بندی ممکن است تغییر کند. در این قسمت از کتابخانه ی یاد گیری ماشین پایتون به نام sklearn استفاده کردیم که تمام کرنل ها و توابع نگاشت را به صورت آماده دارد. ۳ تابع LinearSVC, NuSVC, SVC وظیفه اصلی دسته بندی را در این کتابخانه دارند که ما از SVC که مخفف SVC است استفاده می کنیم.

ورودی های این تابع کرنل مورد استفاده و پارامتر C است. کرنل را گاوسی (rbf) در نظر گرفتیم و پارامتر C را از ورودی می گیریم. با استفاده از کرنل گاوسی، داده ها را به فضای بی نهایت بعدی نگاشت می کنیم و در آن فضا دسته بندی را انجام می دهیم. پارامتر C نیز در واقع میزان اهمیت به اشتباه دسته بندی شدن داده ها و فاصله ی خط جداساز از Support Vectorها را مشخص می کند.

در کنسول می توان با وارد کردن دستور svm test c که پارامتری است که کاربر ورودی می دهد، برچسبهای خروجی دسته بند SVM را برای مجموعه ی تست مشاهده کرد.

```
def svm(self, x_train, y_train, x_test, c_parameter):
   model = SVC(kernel='rbf', C=c_parameter)
   model.fit(x_train, y_train)
   self.svm_classifier = model
   y_pred = model.predict(x_test)
   return y_pred
```

تصویر ۷: دسته بند SVM

Random Forest

الگوریتم جنگل تصادفی یا همان Random Forest یک الگوریتم ترکیبی (ensemble) است که از درختهای تصمیم بهره می گیرد.

الگوریتم درخت تصمیم (decision trees) می تواند به راحتی عملیات طبقه بندی را بر روی داده ها انجام دهد. حال در الگوریتم جنگل تصادفی از چندین درخت تصمیم (برای مثال ۱۰۰ درخت تصمیم) استفاده می شود. در واقع مجموعه ای از درخت های تصمیم با هم یک جنگل را تولید می کنند و این جنگل می تواند تصمیم های بهتری را نسبت به یک درخت اتخاذ نماید.

در الگوریتم جنگل تصادفی به هر کدام از درختها، زیرمجموعهای از داده ها تزریق می شود. برای مثال اگر مجموعه داده دارای ۲۵۵۰ سطر باشد و ۱۰ ستون داشته باشد، الگوریتم جنگل تصادفی به هر کدام از درخت ها ۱۰۰ سطر و ۵ ستون –که به طور تصادفی انتخاب شده اند و زیرمجموعهای از مجموعهی داده ها هست می دهد. این درختها با همین دیتاست زیرمجموعه، می توانند تصمیم بگیرند و مدل دسته بند خود را بسازند.

برای پیاده سازی این بخش، از کتابخانه ی sklearn استفاده می کنیم و از کلاس ensemble که مربوط به الگوریتم های ترکیبی است، تابع RandomForestClassifier را برای دسته بندی مورد استفاده قرار می دهیم و با استفاده از تابع فاز آموزش را بر روی داده های آموزش انجام می دهیم و سپس با استفاده از predict، بر چسب داده های آزمایش را معین می کنیم. با استفاده از دستور random_forrest test می توان بر چسب ها را در کنسول مشاهده کرد.

```
def random_forrest(self, x_train, y_train, x_test):
    model = RandomForestClassifier()
    model.fit(x_train, y_train)
    y_pred = model.predict(x_test)
```

تصویر ۸: تابع دستهبند random_forrest

```
| Josephiny | Thomas | American | The section | The sectio
```

تصویر ۹: برچسبگذاری هر یک از دستهبندها بر روی مجموعه دادهی آزمایش

محاسبهی Validation Set

تابع make_validation_set و make_knn_validation_set برای این قسمت پیاده سازی شده است. هر دوی این توابع عملکرد یکسانی دارند و به ازای هر مستند در فضای برداری یک عدد رندوم با توزیع یکنواخت تولید می کنیم تا با احتمال 0.1 مستند در مجموعه ی Validation و با احتمال 0.9 در مجموعه ی آموزش قرار بگیرد. سپس دو مجموعه را باز می گردانیم.

```
def make_validation_set(self):
    x_train_set = []
    x_validation_set = []
    y_train_set = []
    y_validation_set = []
    for i in range(self.train_size):
        m = random.uniform(0, 1)
        if m < 0.9:
            x_train_set += [self.train_vector_space[i]]
            y_train_set += [self.y_train[i]]
        else:
            x_validation_set += [self.train_vector_space[i]]
            y_validation_set += [self.y_train[i]]
        return x_train_set, y_train_set, x_validation_set, y_validation_set</pre>
```

تصویر ۱۰: تابع ایجاد مجموعهی Validation از مجموعهی آموزش

محاسبهی بهترین پارامتر برای الگوریتم k-NN

در این بخش از میان kهای گفته شده در صورت پروژه -یعنی ۱ و ۵ و ۹- بهترین k را انتخاب می کنیم و در قسمتهای بعدی نیز از آن استفاده می کنیم.

برای این کار از تابع find_best_k استفاده می کنیم. این تابع در صورتی که فلگ print برابر با True باشد، به ازای هر لا از مقادیر گفته شده، تابع knn را با توجه به مجموعه ی آموزش و Validation صدا می زند و سپس با استفاده از خروجی آن و تابع find_metric که در بخش پیاده سازی تابع find_metric توضیح داده شده، معیارهای مختلف را به ازای هر لا در خروجی چاپ می کند. در نهایت بهترین لا را با توجه به بیشترین مقدار accuracy انتخاب کرده و چاپ می کنیم. در صورتی هم که فلگ print برابر False باشد، تنها بهترین مقدار لا محاسبه می شود و چیزی چاپ نمی شود.

```
def find_best_k(self, arr, print_flag):
    max_accuracy = -1
    best_k = None
    x_train_set, y_train_set, x_validation_set, y_validation_set = self.make_knn_validation_set()
    for k in arr:
        y_pred = self.knn(x_train_set, y_train_set, x_validation_set, k)
        f1 = self.find_metric(y_validation_set, y_pred, "f1")
        precision = self.find_metric(y_validation_set, y_pred, "precision")
        recall = self.find_metric(y_validation_set, y_pred, "recall")
        accuracy = self.find_metric(y_validation_set, y_pred, "accuracy")
        if print_flag:
            print("metrics for k = ", k)
            print("f1 score = ", f1)
            print("precision = ", precision)
            print("precision = ", precision)
            print("recall = ", recall)
            print("accuracy = ", accuracy)
            print()
        if accuracy > max_accuracy:
            max_accuracy = accuracy
            best_k = k
        if print_flag:
            print("best k is: ", best_k)
            print("best accuracy score is: ", max_accuracy)
        return best_k
```

تصویر ۱۱: پیاده سازی تابع یافتن بهترین پارامتر k

در کنسول می توان با دستور k = 9 خروجی این تابع را مشاهده کرد. مشاهده می شود که k = 9 بهترین k بوده است.

تصویر ۱۲: خروجی دستور best k

محاسبهی بهترین پارامتر برای الگوریتم SVM

برای پیاده سازی این بخش از تابع $\operatorname{find_best_c}$ استفاده شده که عینا مشابه تابع $\operatorname{find_best_k}$ عمل می کند که در بخش قبل توضیح داده شد. با این تفاوت که به جای تابع sym تابع sym صدا زده می شود و مقدار بهترین C از بین مقادیر C . 1. 5.1 و C انتخاب می شود.

```
def find_best_c(self, arr, print_flag):
    max_accuracy = -1
    best_c = None
    x_train_set, y_train_set, x_validation_set, y_validation_set = self.make_validation_set()
    for c in arr:
        y_pred = self.svm(x_train_set, y_train_set, x_validation_set, c)
        f1 = self.find_metric(y_validation_set, y_pred, "f1")
        precision = self.find_metric(y_validation_set, y_pred, "precision")
        recall = self.find_metric(y_validation_set, y_pred, "recall")
        accuracy = self.find_metric(y_validation_set, y_pred, "accuracy")
        if print_flag:
            print("metrics for c = ", c)
            print("f1 score = ", f1)
            print("precision = ", precision)
            print("recall = ", recall)
            print("accuracy = ", accuracy)
            print("accuracy = accuracy)
            best_c = c

if print_flag:
        print("best c is: ", best_c)
        print("best accuracy score is: ", max_accuracy)
        return best_c
```

تصویر ۱۳: تابع یافتن بهترین پارامتر C

با دستور c=1.5 می توان خروجی مورد نظر را مشاهده کرد که نشان می دهد c=1.5 بهترین مقدار بوده است.

```
hest c
metrics for c = 0.5
f1 score = 0.6254545454545455
precision = 0.5308641975308642
recall = 0.7610619469026548
accuracy = 0.5928853754940712
f1 score = 0.5738396624472574
precision = 0.5483870967741935
recall = 0.6017699115044248
accuracy = 0.6007905138339921
metrics for c = 1.5
f1 score = 0.5907172995780591
precision = 0.5645161290322581
recall = 0.6194690265486725
accuracy = 0.616600790513834
metrics for c = 2
f1 score = 0.5822784810126581
precision = 0.5564516129032258
recall = 0.6106194690265486
accuracy = 0.6086956521739131
best c is: 1.5
best accuracy score is: 0.616600790513834
تصویر ۱۴: خروجی دستور ۱۴
```

۱۳

بهبود سیستم بازیابی فاز اول پروژه

در این قسمت، ابتدا با استفاده از دادههای آموزشی هر چهار دسته بند را آموزش می دهیم. سپس دادههای فاز ۱ را که در این قسمت، ابتدا با استفاده از هر چهار دسته بند بر چسب می زنیم و ذخیره می کنیم. این ذخیره سازی به کمک دستورات تصویر زیر انجام می شود.

```
naive_bayes phase1
y's saved for naive_bayes method
knn phase1
y's saved for knn method
svm phase1
y's saved for svm method
random_forrest phase1
y's saved for random_forrest method
calculate majority vote
y's saved for majority_vote_prediction method
do not be not be
```

دستور calculate majority vote رای های مستندات را مجتمع می کند و بهترین بر چسب را انتخاب می کند. نحوه ی کار آن در بخش تابع اجماع نظر دسته بندها برای تعیین بر چسب هر مستند توضیح داده می شود. سپس این بر چسبها را در فایلی ذخیره می کند.

حال می توانیم با توجه به ورودی کاربر که یا ۱ یا ۱- است، مستندات مرتبط با آن دسته را در خروجی به کاربر نمایش دهیم. به این صورت که پس از فیلتر کردن مستندات فاز اول بر اساس دسته، مطابق همان چیزی که در فاز اول پیاده سازی شده است، مستندات را بر اساس معیار ltc-lnc مرتب سازی می کنیم و به کاربر نمایش می دهیم. در ادامه توابع این بخش را توضیح می دهیم.

€ classifier.py	10.7 kB	14:11
data	4 items	Yesterday
🙋 ir_system.py	46.4 kB	14:29
knn_y_prediction	8.8 kB	14:21
main.py	19.3 kB	13:57
majority_vote_prediction_y_prediction	10.5 kB	14:24
anaive_bayes_y_prediction	8.7 kB	14:18
hase1_classifier_data	470.2 MB	13:30
phase2_command	481 bytes	14:16
pycache	2 items	14:29
andom_forrest_y_prediction	20.6 kB	14:24
svm_y_prediction	20.6 kB	14:24
test_classifier_data	248.1 MB	13:30

تصویر ۱۶: فایل های جانبی ذخیره شده که بر چسبهای هر دستهبندی را نگه میدارند.

تابع اجماع نظر دسته بندها برای تعیین برچسب هر مستند

در این قسمت، ابتدا برچسبهای ذخیرهسازی شده برای هر دسته بند را بارگذاری می کنیم و سپس برای هر مستند نظر اجماع دسته بندها را به عنوان برچسب آن مستند در نظر می گیریم.

نکتهای که باید به آن توجه کرد این است که چهار دسته بند داریم و لذا ممکن است دو دسته بند نظر به برچسب ۱ داشته باشند و دو تای دیگر نظر به برچسب ۱-. در این حالت برچسبی که دسته بند Random Forest برای مستند انتخاب کرده است به عنوان نتیجه ی نهایی در نظر گرفته می شود، زیرا با آن به دقت بالاتری می رسیدیم. همان طور که گفته شد، این اجماع نظر به کمک دستور calculate majority vote انجام می گیرد.

جستجو بر اساس دسته

تابع این بخش، phase1_query است که در IRSystem طراحی شده است. ابتدا نظر جمعی دسته بندها را لود می کنیم و از آن برای بر چسب زدن هر مستند فاز اول استفاده می کنیم. سپس مانند فاز قبلی پروژه عمل می کنیم، با این تفاوت که جستجو را در دسته ی مشخص شده انجام می دهیم.

در این قسمت برای جستجو با انتخاب دسته ی مورد نظر دستور phase1 query در کنسول زده می شود. سپس از کاربر می خواهیم دسته ی مورد نظر را انتخاب کند و بعد از آن مستندات را با توجه به ورودی کاربر فیلتر می کنیم. پس از آن کاربر می بایست پرسمان خود را وارد کند و پس از تصحیح پرسمان و رتبه بندی مطابق فاز قبل، نتایج به کاربر نشان داده می شود.

نمونه هایی از اجرای این بخش در ادامه دیده می شود.

```
propore english
creation was successful
create positional english
creation was successful
phase1 guery
Please Enter Zone Of Search(1 > Most View ,-1 > Less View): 1
Enter Your Query: diamac chem
no spell correction needed!
document 2054: ['slatheia*, 'not', 'normal', 'age', 'cure'], ['than', 'million', 'peopl', 'wordwid', 'suffer', 'alzheia*, 'diseas', 'number', 'expect', 'increas', 'drastic', 'come', 'year', 'no', 'real', 'progress',
ttc-lnc score: 0.20062230780457598
document 2051: [['get', 'back', 'work', 'after', 'career', 'break'], ['if', 'youv', 'taken', 'career', 'break', 'nom', 'look', 'return', 'workfore', 'would', 'consid', 'take', 'internship', 'career', 'reentri', 'expectic-lnc score: 0.1809856065322041
document 1051: [['sugt', 'back', 'work', 'after', 'career', 'break'], 'do', 'argu', 'outreason', 'oppon', 'prove', 'them', 'wrong', 'most', 'all', 'win', 'right', 'philosoph', 'daniel', 'h', 'cohen', 'show', 'most', 'common', 'form', 'argument', 'w
ttc-lnc score: 0.16959085780750133
phase1 query
Please Enter Zone Of Search(1 > Most View ,-1 > Less View): -1
Enter Your Query: diamac cohen
no spell correction needed!
document 795: [['tough', 'truth', 'plastic', 'pollut'], ['artist', 'diama', 'cohen', 'share', 'some', 'tough', 'truth', 'plastic', 'pollut', 'ocean', 'live', 'some', 'thought', 'free', 'ourselv', 'plastic', 'gyre']]
ttc-lnc score: 0.475442886582986
```

تصویر ۱۷: اجرای دستو ر phase1 query برای هر دو دسته با پر سمان ۱۲: اجرای دستو ر

```
phaseI query

Please Enter Zone Of Search(1 -> Most View ,-1 -> Less View): 1

Enter Your Query: pericillin
no spell correction needed!

document 2004: [['do', 'do', 'when', 'antibiot', 'dont', 'work', 'ani'], ['penicillin', 'chang', 'everyth', 'infect', 'had', 'previous', 'kill', 'were', 'sudden', 'quick', 'curabl', 'yet', 'maryn', 'mckenna', 'share',
ltc-lnc score: 0.5091126097043745

phaseI query

Please Enter Zone Of Search(1 -> Most View ,-1 -> Less View): -1

Enter Your Query: penicillin
no spell correction needed!

Query Not Found in Your Zone
```

تصویر ۱۸: اجرای دستور phase1 query برای هر دو دسته با پرسمان

ارزیابی نهایی

در این بخش ابتدا تابع find_metric را پیاده سازی می کنیم که با توجه به رشته ی ورودی که همان معیار مورد نظر است، مقدار آن معیار را با توجه به دو ورودی دیگر که دو لیست از مقادیر واقعی و مقادیر پیش بینی شده هستند، خروجی می دهد. سپس برای هر کدام از دسته بندها معیارهای خواسته شده را محاسبه می کنیم که در ادامه توضیح داده می شوند.

پیادهسازی تابع find_metric

تابع find_metric سه ورودي به نامهاي y_pred ،y_test دارد كه:

- ورودی y_test شامل گروههای واقعی مستندات آزمایش (test.csv) است.
- ورودی y_pred شامل گروههای پیش بینی شده ی مستندات آزمایش است.
- ورودی metric نیز معیار ارزیابی را مشخص می کند که یک رشته است و مقدار آن برابر با یکی از مقادیر precision ،accuracy ،f1

نحوه ی عملکرد این تابع نیز به این صورت است که با توجه به ورودی ها تعداد True Negative، True Positive، معیار خواسته شده False Positive و False Negative را به دست می آورد. سپس با توجه به فرمول های اسلایدها، معیار خواسته شده را خروجی می دهد.

Naïve Bayes معیارهای ارزیابی برای

در این قسمت، با وارد کردن دستور metric] naïve_bayes] در کنسول، می توان مقدار معیار ورودی را برای دسته بند Naïve-Bayes مشاهده کرد.

accuracy naive_bayes

accuracy is: 0.6509803921568628

precision naive_bayes

precision is: 0.6159420289855072

recall naive_bayes

recall is: 0.7024793388429752

f1 naive_bayes

f1 is: 0.6563706563706564

تصویر ۱۹: مقدار معیارهای ارزیابی برای دسته بند Naïve Bayes

معیارهای ارزیابی برای k-NN

با وارد کردن دستور k metric k مشاهده کرد. با توجه به اورد کردن دستور k مشاهده کرد. با توجه به این که در این بخش معیارها برای بهترین k خواسته شده بود، k را برابر با k قرار میدهیم.

accuracy knn 9
accuracy is: 0.5647058823529412
precision knn 9
precision is: 0.5384615384615384
recall knn 9
recall is: 0.5785123966942148
f1 knn 9
f1 is: 0.5577689243027888

(k = 9) k-NN تصویر ۲۰: مقدار معیارهای ارزیابی برای دستهبند

معیارهای ارزیابی برای SVM

دستور metric] svm c را در کنسول وارد می کنیم که با توجه به صورت پروژه، مقدار c را برابر با بهترین C، یعنی 1.5 قرار می دهیم.

معیارهای ارزیابی برای Random Forest

از دستور metric] random_forrest استفاده می کنیم که مقدار معیار ورودی را تحت دسته بندی metric] استفاده می کنیم

accuracy random_forrest

accuracy is: 0.6431372549019608

precision random_forrest

precision is: 0.6728971962616822

recall random_forrest

recall is: 0.628099173553719

f1 random_forrest

f1 is: 0.6419753086419753

تصویر ۲۲: مقدار معیارهای ارزیابی برای دسته بند Random Forest

نتیجه گیری نهایی معیارهای ارزیابی

با توجه به تصاویر ۱۹، ۲۰، ۲۱ و ۲۲ می توان نتیجه گرفت که:

- دسته بند SVM بیشترین accuracy را داشته است.
- دسته بند Random Forest بیشترین مقدار precision را حاصل کرده است.
 - دستهبند Naïve Bayes بیشترین مقدار F1 و recall را داشتهاست.
- معیارهای ارزیابی دسته بند k-NN نسبت به سایر دسته بندها پایین تر بوده است.

نحوهى تقسيم وظايف

وظایف اختصاص یافته به هر فرد به شرح زیر بود:

نيما جمالي

۱- پیش پر دازش دادهها و پیادهسازی روش ntn

۲- ساخت دو مدل فضای برداری و توابع مربوط به آنها

تا پیاده سازی دسته بند k-NN و توابع جانبی آن k-NN

Validation و ساختن مجموعهی k پیاده سازی تابع محاسبه ی بهترین پارامتر k

سپهر فعلی

۱- طراحی دستهبند Naïve-Bayes و توابع مربوط به آن

۲- طراحی توابع محاسبهی معیارهای ارزیابی

۳- مقایسهی معیارهای ارزیابی برای دستهبندها

سينا كاظمى

۱- طراحی دسته بندهای SVM و Random Forest

٣- طراحي تابع اجماع نظر دستهبندها

۴- طراحی تابع پاسخ پرسمان، بسته به دستهی ورودی



- [1] https://scikit-learn.org/stable/modules/svm.html
- [2] https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html