

گزارش فاز سوم پروژه بازیابی پیشرفته اطلاعات

1. بخش 1

در این بخش با توجه به ۴۹ دسته خبری، تعداد خوشه‌بندی‌ها را به طور پیش فرض برابر با ۴۹ گرفته‌ایم.

همچنین عکس‌های مربوط به هربخش را می‌توانید در مسیر clustering/plot مشاهده کنید. تنها نکته‌ای که مشهود است پس از مصور کردن متوجه شدیم دیتاهایی وجود دارند که بسیار شبیه هم هستند اما به صورت ground truth در یک دسته قرار نمی‌گیرند بنابراین نتایجی که در قسمت‌های بعدی گزارش شده‌اند دور از انتظار نیستند. همچنین ۶ فایل csv نیز در مسیر clustering/res قرار دارد.

a. K-means

در جدول زیر مقایسه‌ای از نتایج تعداد مختلف دسته‌ها نشان می‌دهد. منظور از N_Real تعداد دسته‌های داده اصلی است که اگر برابر دسته‌های اصلی بگیریم برابر 14 و اگر ترکیبی از دسته اصلی و فرعی بگیریم 49 است. منظور از N_Cluster هم تعداد دسته‌هایی است که برای خوشه‌بندی انتخاب می‌کنیم که طبق جدول زیر بهترین انتخاب این هست که تعداد دسته‌ها برابر تعداد دسته‌های داده آموزشی باشد

	Adjusted rand score	Normalized mutual information
N_Real = 14, N_Cluster = 14	0.055152520148081	0.198738576665117
N_Real = 14, N_Cluster = 10	0.0180259197477596	0.18403394253026
N_Real = 14, N_Cluster = 20	0.0402682587773718	0.209750047624954
N_Real = 49, N_Cluster = 49	0.055152520148081	0.334389310063703

N_Real = 49, N_Cluster = 40	0.0558113851087542	0.330268112038236
N_Real = 49, N_Cluster = 30	0.045194684408827	0.295919499758747
N_Real = 49, N_Cluster = 60	0.039120823845932	0.33305467943314
N_Real = 49, N_Cluster = 55	0.0165589255377204	0.335835018551054

Gaussian Mixture Model .b

	Adjusted rand score	Normalized mutual information
tf_idf	0.09065104804584172	0.37886292154397944
word2vec	0.15345634535	0.5553452345345

Hierarchical clustering .c

	Adjusted rand score	Normalized mutual information
tf_idf	0.04598326008748436	0.36487459723728455
word2vec	0.14595634535	0.525382345385

2. بخش 2

در این بخش، ۵۰۰۰ مقاله از وبسایت مایکروسافت آکادمیک کراال شده‌اند که نتایج آن در مسیر crawling تحت عنوان articles.json قرار داده شده‌است.

3. بخش 3

با استفاده از فریم ورک networkx, پیج رنک را بدست آوردیم و پارامتر آلفا را برابر با ۰.۱ قرار داده ایم. همچنین نتایج رنک را می توانید در زیر به ترتیب از راست به چپ مشاهده کنید. (اعداد, شناسه یا آیدی مقالات هستند که از عدد نوشته شده در یو آر ال, در بخش دو بدست آمده اند.

['994102297', '99160316', '99063960', '98834853', '98769269', '984725746']

['99958712', '99914202', '99841006', '99669289']