

知能情報実験 III（データマイニング班）YouTube 急上昇 に乗る動画のタグから、現在の流行を掴み、顧客をアシ ストする

グループの学籍番号列挙 215708A, 215758G, 215752H,215764B

提出日：2023 年 8 月 3 日

概要

本文書は、知能情報実験班 II(I データマイニング) 班 GP5 班の報告書である。GP5 班は、「YouTube 急上昇に載る動画のタグから、現在の流行を掴み、顧客をアシストする」ことを目的とし、現在の 流行となる動画を、急上昇に乗っている動画のタグから掴むことができる機械学習モデルの開発を 行っている。データセットは、YoutubeAPI を用いて情報を取得し.csv ファイル形式で構築して いる。学習モデルはクラスタリングを使用した。しかし、結論として、このプログラミングは成功しなかった。原因はクラスタリングの使用はあまり適切でないことと、データ量が不足していたことである。本来、機械学習によってクラスタリングする時には単語を数値化させる必要があるが、今回の我々はタグの出現回数を単語の数値として扱ったためであると考察する。

目次

1	実験の目的と達成目標	2
2	テーマ：急上昇とは	2
3	実験方法	2
3.1	データセット構築	2
3.2	モデル選定	3
3.3	パラメータ調整	3
4	実験結果	3
5	考察	8
6	意図していた実験計画との違い	8
7	まとめ	9

1 実験の目的と達成目標

私達 GP5 班の実験目的は「YouTube 急上昇に乗る動画のタグから、現在の流行を掴み、顧客をアシストする」ということである。YouTube の API から急上昇に乗っている動画のタグの情報を取得し、機械学習を使って現在の流行であるキーワードを抽出する。このキーワードは、今現在注目されている、もしくはされ始めているものであると考えられる為、これを新人動画クリエイターや企業などに提供することで、現在の流行のキーワードを共有し、アシストすることを目標とする。

2 テーマ：急上昇とは

YouTube でバズると言われると思ひ浮かぶ事として急上昇ランキングであろうと思われる。急上昇とは [1] によると現在話題になっていることであり、視聴者が動画をコメント、評価することで視聴回数の伸びに寄与する。本グループでは急上昇における動画のタグから、今トレンドとなっているキーワードを抽出できるのではと考え、対象問題として設定した。

3 実験方法

今回の実験では、トレンドとなっているキーワードは急上昇動画の中のタグでも頻出しているだろうと予測し、1 週間ごとに急上昇に乗った動画のタグを全てまとめて特徴量を抽出することで、その週に流行したタグを抽出した。また、トレンドとなっている内容であれば、何週間かにわたって急上昇に乗り続ける可能性が高いと予測し、過去の 1 週間の急上昇データと現在の急上昇データを比較することで、より大きな流行となっているキーワードの抽出を行いたいと考えた。今回の実験では下記の三段階の手順を踏む。

1. 自分の Google アカウントから「Youtube Data API」を生成しそれを用いて、現在の急上昇に乗っている動画の情報を取得する。
2. 取得した動画から頻度の多いタグ順に並べて csv ファイルに保存する。
3. 先ほどの csv ファイルは tag の数が多く出るため frequency が 2 以上であるものだけにフィルターをかけて新たに csv ファイルとして保存する。
4. 2 と 3 を数回繰り返し 1 週間ごとのデータを抽出する。
5. 過去と現在のデータを読み込み機械学習 (クラスタリング) にかける

3.1 データセット構築

今回は、既存のデータセットを用いるのではなく、[2] による YoutubeAPI を用いて情報を取得することで、新たにデータセットを作成することにした。取得する情報としては、現在の YouTube の急上昇 50 位までのチャンネル名、動画のタイトル、動画につけられている tag、投稿日、チャン

ネル ID、視聴回数、高評価である。取得した情報は csv ファイルに変換して保存している。この csv ファイルから tag の頻出度を数え、頻出度が 2 以上であるものをトレンドであると定義し新しく csv ファイルで保存している。

3.2 モデル選定

tag それぞれの重みをつけてジャンル分けをするために k-means のクラスタリングを使用することにした。なぜなら、私たちが今回行う実験では YouTube の急上昇から tag の情報を抽出するだけだと統計計算だけで済むから機械学習を取り入れるためクラスタリングを取り入れた。クラスタリングを行うことで視覚化しユーザにあったジャンルをアシストできると考えたからである。

3.3 パラメータ調整

クラスタ数を 4、初期セントロイドを 42 とし、その他はデフォルト値のままにしている。

4 実験結果

今回の実験では YouTube API を用いて YouTube における急上昇動画 50 本の情報を抜き出した (図 1)。そこからタグとして使用されている単語を抜き出し、出現回数と共に csv ファイルに記録した (図 2)。また、その際に比較的出現回数の多いタグを厳選したファイルも作成した (図 3)。その後クラスタリングによる機械学習を行ったが、実行した可視化するためのグラフが予想したものとは異なっていた。

youtubetop_50

	0	channelTitle	publishedAt	channelId	title
0	E8BgYs7hUJU	ジャにのちゃんねる	2023-08-02T10:00:19Z	UC2ahDZWkakOTxCoF-uMAg	#259 【ふっか博士召喚】 横山田の同席。
1	-Bnk_Tjcm3Y	HikakinTV	2023-08-01T10:00:13Z	UCZF_ehICEBPop_sldpBUQ	20億円のはきん新居ハウスツアー！ 超巨大室内連水プール&庭付きの家
2	JGDyofr7a3U	LOVECARSTV	2023-08-02T10:15:01Z	UCtLo4mwb3ObCDZ4m8b8u7fA	トヨタ ランドクルーザー 250 実車詳細レビュー！ コレ欲しい！ 魅力的！台が整備！
3	QdQWdStcMAE	スカイピース	2023-08-02T11:00:09Z	UC8_wmm5DX9mb4jrLiw8ZYzw	【カルマ帰郷】 牛角で人気ランキング1位→10位まで下りたよ10！！！！
4	ZkxZfHors-og	JFATV	2023-08-02T11:13:216Z	UCgleUSV91-FfmCayG4S8Cw	[LIVE] 第103回天童杯 F C 町田ゼルビア vs.アルビレックス新潟 ラウンド16 (48戦)
5	FHsehL2QaTfw	舞元香村	2023-08-01T18:56:34Z	UCJubNhcCofXa8wntqbwLg	【#にじ生2023】 熱狂！ にじさんじ 甲子園2023 vol.9 【にじさんじ/舞元香村】
6	6x0xR5y4F50	カラフルピーチ	2023-08-01T08:00:07Z	UC07b7ThbAkAro6t6jgrt8B1HhA	成人にエロがきたみたい？ 【ミニゲーム】
7	1vYXGRyHyo0	トヨタYouTubeショールーム	2023-08-02T01:17:43Z	UCvxtexS8-gjyebtbQd4lg	トヨタ自動車 新車ランドクルーザー ワールドプレミア
8	5OHp92FO3	KITERU イキタール	2023-08-01T11:01:28Z	UC8BLyb67OxGtnaGuJArlqiw	【人生初】 韓国高校生が初めて日本に来て衝撃!!! 教科書に書いてあった国と全く違って1日目の24時間が驚きの連続！
9	x-eFEnQ3QLU	FNNプライムオンライン	2023-08-01T11:30:31Z	UCoQ8JMcwcmXRShBFAt5lw	長男・玄一前副社長に大失態！さつ ビッグモーター幹部ら「電流警備点検」被検入手
10	uKfEnOGGMl	jun channel	2023-08-01T15:04:06Z	UCx1nAvrVtDaaGmCM5e8ofuq	夏のホラゲ配信-1【Shadow Corridor dlc】
11	21EP0H+2Es	不破 遼 / Fuwa Minato 【にじさんじ】	2023-08-01T13:02:36Z	UC6wvdADTJ88OfbYjPaADA	【重大発表アリ】 FuwaHiba Session LIVE 【不破遼/逢会響彦/にじさんじ】
12	SKyQ3Tdf_E	TREASURE (트레저)	2023-08-01T15:00:05Z	UCxdhXYOCuUyWprEqa4ZQdHA	TREASURE - 'BONA BONA' DANCE PRACTICE VIDEO
13	P77PjgKMDMU	【東海テレビ公式】 ドラHOTpress	2023-08-01T12:58:30Z	UC2G0uq2KGfKonnG8yWeNZ_Q	8/1 中井x飯神 ハイライト
14	AW43GJmsIM	バカラジオ / Bankaradio	2023-08-01T09:05:27Z	UCT5bVnm5mgD8xdoZ8FhA	ボケカの転売で荒療ぎする小学生
15	ngp5o5X0XQ	SAWYAN CHANNEL / サワヤン チャンネル	2023-08-01T11:00:41Z	UCdndn3NwRfouyCAVnRTzPQ9g	【番組情報】 減薪中のヤンと1日全く同じ食生活したらまさかの結果に、、、
16	lyf8-Y6A9sk	びへチャンネル	2023-08-01T11:00:38Z	UCdHwR1u4fR8gbdhnd5GWAQ	ボケモン界の歌姫メロエッタを歌いたい【ボケモンSV】 【ゆっくり実況】
17	-L8GzT3CK0	アニメ「おでかけ子メ」チャンネル	2023-08-01T11:00:09Z	UCjfn8eLsvEGyK_b5YU51KQD	アニメ「おでかけ子メ」第1話【映画】
18	8LGZDFrktE	日テレNEWS	2023-08-01T00:26:10Z	UCuTAXTewrhettOx3zgsakJBQ	【ビッグモーター不正】「新たな疑念」 現役従業員「整備せず納車」証言 会社側の回答は…
19	ObIRwBp8b-2Q	JYP Entertainment	2023-07-31T08:58:09Z	UCaOETy8CBUSjtt6zhTrZgg	ITZY 'CAKE' M/V @ITZY
20	nVPy7hpkuzw	MBS NEWS	2023-07-31T09:33:16Z	UCaKJhKQ73xf1Pk5hR8Z2GQ	【ビッグモーター】 副社長のバウハラLINE「死刑裁判死刑…」 「実質支配の重責親子との関係断つための民事再生生活」 山岸久前弁護士の見解【MBSニュース解説】 (2023年7月31日)
21	jhDMJN10xkY	HikakinGames	2023-08-01T11:00:11Z	UCX1xppLuvQ3bUcJ8ajyA	10年間ありがとうごさしました。
22	v7ePavkhFH0	サントリー公式チャンネル (SUNTORY)	2023-07-31T12:24:06Z	UCyHbnB7_27hyaYgCf9Z4g	【プロ野球界のレジェンドが集結】 サントリードリームマッチ2023アーカイブ配信
23	CY5aYBtLs	にこちゃん放送局	2023-07-31T11:30:03Z	UCmWkowGJ177P9jWWz85vTw	【裏面密着】 デイズニールとピクサー映画の声を聴くもこれになりました
24	aRWna815ys	THE RAMPAGE from EXILE TRIBE	2023-07-31T12:00:07Z	UCY3GYW6aVaxDwZvW86wKQ	THE RAMPAGE / Summer Riot ~ 祭典夜 ~ (MUSIC VIDEO)
25	3WS5AMG3yJ4	Apex Legends	2023-07-31T15:00:11Z	UCQZV8MZTHAS1Q79wVWUC3A	Apex Legends: Resurrection Launch Trailer K8 Code - Part 2
26	6-A6eNzJus	匿名実業 / Shine Yuika	2023-07-31T12:42:57Z	UC_4Xjgacpqs5Jc05ncncpQ	【#にじ生2023】 #9 にじさんじ高校3年生〜【にじさんじ/匿名実業】
27	K_f1aJbPxa0	斎藤雄	2023-07-31T00:02:37Z	UCeK1fLpJtHLYZUc0rVfK4cw	超RIZIN2の感想
28	JOpWVU_F13c	RKB毎日放送NEWS	2023-07-31T09:11:54Z	UCgrmh4AbzN-QzMMgJ2zNf0g	ビッグモーター 現役店長と社員が証言 「LINEで1日2000件の業務連絡」「第1本生えていたら減点」
29	14MqKJdIR30	oricon	2023-07-30T11:33:35Z	UCzZwK2uAgR0vsa4FytsCQ	【超RIZIN2】 前後未来、クラフトに豪華の一歩負け 試合直後の心機明かす「強かったです…」 「超RIZIN 2」 試合後インタビュー
30	7hT7f0kwRA	ChroNoR	2023-07-31T10:00:12Z	UCz6vmbjgqfT9xUcD8Bp65Q	【このツイ特ごうタイズ】 にじさんじライバーのバズったツイートを挙げて1！くるなん
31	YM_3DNXB41c	らっだぁ	2023-07-31T09:03:42Z	UCiNnVSK1uy2zHfTaXZYXopw	売付けたら自分以外関係者でワンオペ状態になっていたらだぁ【VCRグラセフ/VRGTA】
32	Fickd9T2BTQ	ONE PIECE公式YouTubeチャンネル	2023-07-30T01:00:04Z	UCdAHaWckdptT5XhN2Xe6BUQ	ONE PIECE 1071 話予告「ルフィの最悪地点 到達！」「オデス。」
33	oGfTKGwnKxM	ホモサビ	2023-07-30T13:28:44Z	UCd0hscDvuzrRo6Rk7JPQMA	カメムシ漬してレモンを絞るとコウラの香りになる
34	Uk6dF4rm2w	おちちんゆー	2023-07-30T08:28:45Z	UCyprf0dH4H-zuL7mnpP-JK2yQ	画にセミ響けかけたwww
35	zEjgpgysgnc	ピロギとニコシーちゃんねる	2023-07-30T15:20:23Z	UCs49FF01aYShCEanCzpgkQ	【超RIZIN 2】 試合を録えて【ニコヤスP】夢劇場
36	Kz7H7Y859p	△△△と仲良し ZOOM	2023-07-30T09:00:10Z	UCM3jwN8R5-WBjvqT5DuaEQ	【9月1日 可也4K】 声優 Seven (feat. Luffj [Song Kook FanCam] @SBS Inkigayo 230730
37	OUJASAna3ck	東京03 Official YouTube Channel	2023-07-30T03:00:07Z	UCconics8Q0z5vYTamGtq-Lsw	東京03「スィッチ」/「第23回東京03単独公演「ヤな返風」」より
38	I-PnC1YrGM	Snow Man	2023-07-28T14:00:09Z	UCuFPaemAmMR8R5cHqyZ3dQ	Snow Man "Dangerholio" Music Video YouTube Ver.
39	Tv4SVnqJdg	RIZIN FIGHTING FEDERATION	2023-07-30T14:08:11Z	UCZZZOUgJWdRdM8_5sqtpYwQ	【勝者と敗者】 試合直後の選手の見解に密着【超RIZIN2 / RIZIN】
40	SRkOt1_uQE	King & Prince	2023-07-30T03:00:13Z	UCSwwcOnzASKE6vofvBn6ATA	King & Prince「なにも」の2 @「音楽の日2023」
41	ocDLPLKx67o	BANGTANTV	2023-07-29T13:00:08Z	UCLAKepWjyImX3toFvFvYQ	【6月8日】EP.15 SUGA with 登壇
42	Yq7jrm41Fo	倉藤正徳の倉ちゃんTV	2023-07-30T12:03:27Z	UCQ7oQLDXVNoqb2gdJ3s-mw	【速報】 社説決着！ 彼らこそ格闘技の醍醐味【超RIZIN2】
43	L_gm5o5QRE	はじめてのしゃ〜り (hajime)	2023-07-30T10:55:54Z	UCd4P6P9R9Kv7wChYUwqzw	【ボケカ】 最後の決闘を100万倍楽しめてみた！！！！
44	uc4yQOQhR90	saotnik	2023-07-30T11:00:40Z	UC79849JZLb4p4H4pGzQzU3nA	予備にない結果を返した山本響良663円高回復作【スト・騎GTA】
45	NkYfVfYz4Q	タイピ-日記/taipi	2023-07-30T11:00:13Z	UCDDH4F4w8P9bPjP9kGQ2A	全身ノミのうらなちだろ子猫をシャンプーしてあたろ。。
46	gm5ShyA236KM	スーツ 交通 / Suit Train	2023-07-30T11:00:08Z	UCdBR2bnFAavDfHqHqQ7a9Q	珍しく東海道新幹線が遅れているの乗れません！
47	-8OFhlgBzmE	スカイピース	2023-07-31T11:00:12Z	UC8_wmm5DX9mb4jrLiw8ZYzw	【ありえない】横浜アリーナのライブで問題が発生しすぎて3150
48	uo33dg7BtMo	ペス	2023-07-30T11:00:43Z	UC0YsAIJUC0yx1cCbshABidA	【仲直りおしゃべり】 カラフル友が集まったら盛り上がりすぎて3150
49	AqEzeiqsKEI	Marine Ch. 宝鐘マリン	2023-07-30T13:03:03Z	UCcZJfH08K0Vv4wQG1vklJvg	【#宝鐘マリン生誕LIVE2023】 3rd Generation Girls Band IIIII【本ロライブ/宝鐘マリン】

図1 今急上昇の動画 50 本のデータ

w_f_2023-08-03

keyword	frequency
MLB	1
メジャーリーグ	1
EPL	1
プレミアリーグ	1
SPFL	1
スコティッシュプレミアシップ	1
三笥薫	1
セリエA	1
エンゼルス	1
大谷翔平	1
ダルビッシュ有	1
吉田正尚	1
富安健洋	1
FAカップ	1
古橋亨梧	1
YG Entertainment	1
YG	1
와이지	1
K-pop	1
트레저	1
트레저	1
TREASURE	1
CHOI HYUN SUK	1
JIHOON	1
YOSHI	1
JUNKYU	1
MASHIHO	1
YOON JAE HYUK	1
ASAHI	1
BANG YE DAM	1
DOYOUNG	1
HARUTO	1
PARK JEONG WOO	1
SO JUNG HWAN	1
최현석	1

図2 急上昇の動画50本のデータからタグの出現数をまとめたデータ

w_f_2023-08-03_filt

keyword	frequency
マネーの虎	2
岩井良明	2
マネーの虎 岩井	2
マネーの虎 現在	2
投資	2
融資	2
資金調達	2
起業	2
ビジネス	2
経営	2
UCTyKZzmKi95wxmCg9rU-j6Q	2
モノリスジャパン	2
れいわのとら	2
TikTok	2
アイルトンモカ	4
本当は不良なのに	2
神回	2
お笑い	2
面白い	2
ゲーム実況	3
ヒカキン	2
ひかきん	2
ユーチューバー	2
YouTuber	2
Vtuber	4
葛葉	3
くずは	2
kuzuha	2
吸血鬼	2
にじさんじ	3
nijisanji	3
いちから株式会社	2
ichikara	2
anime	2
アニメ	2

図3 図2において2回以上出現したタグをまとめたデータ

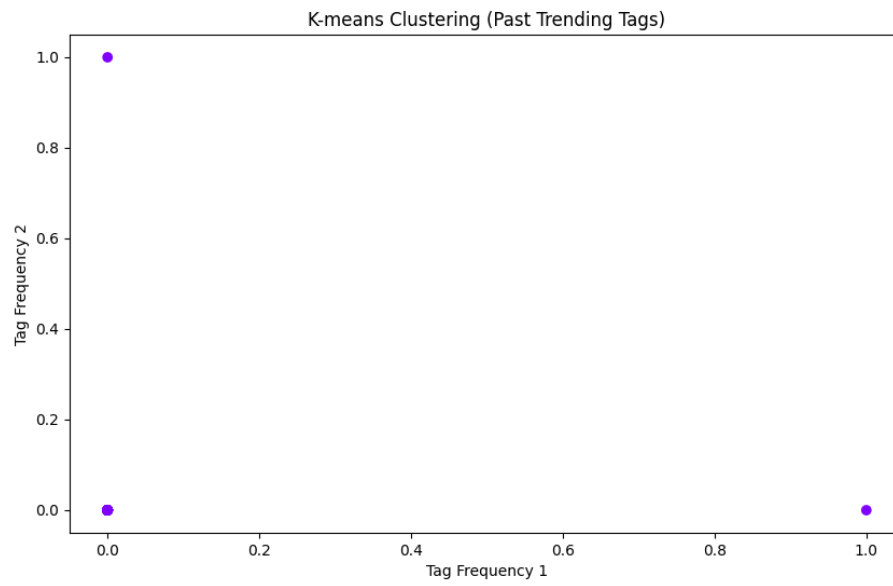


図 4 過去のデータのグラフ

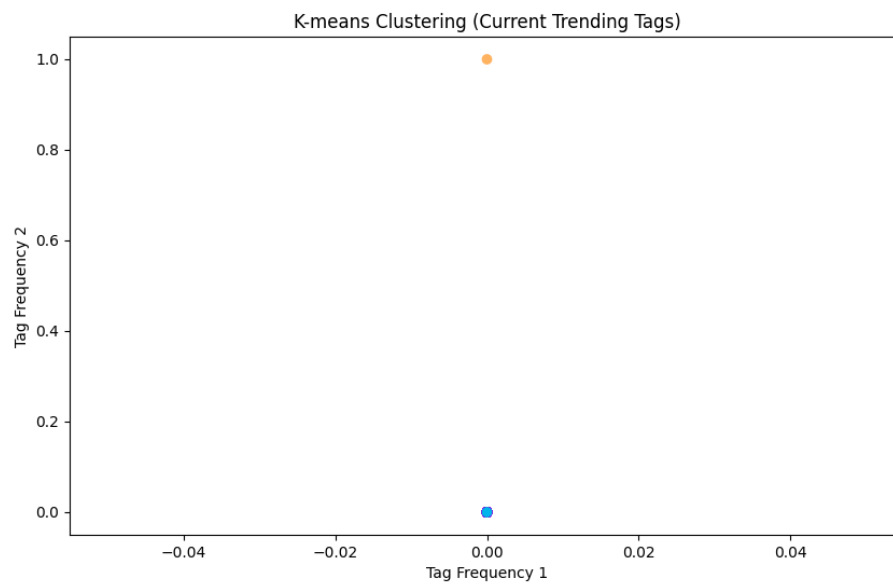


図 5 今週のデータのグラフ

5 考察

今の実験で採用した方法では、目的を達成できなかった。理由としては、以下の2つが考えられる。

1. 機械学習の学習方法が今回の目的に適していなかった。

今回使用した機械学習は主に、クラスタリングの k-means 法と CountVectorizer である。本来、機械学習によってクラスタリングする時には単語を数値化させる必要があるが、今回の我々はタグの出現回数を単語の数値として扱ったため上記ような結果になってしまったと考える。また、CountVectorizer を用いた機械学習も並行して行ったが、そもそも CountVectorizer ではあるテキストから出現する単語のカウントを特徴量にする手法であるが、今回の実験では単語と単語同士で出現する CountVectorizer を使用してしまったため、出現した単語数がカウントしたベクトル化できなかったと考える。

2. データの量が不足していた。

急上昇というジャンルは比較対象が 50 しかなく、機械学習を行うにはあまりにも少ないデータであった。

6 意図していた実験計画との違い

今回の結果を経て、反省点が多く挙がり、大きく感じた反省は、「目標の共有」である。一回目の授業では以下の目標ではなく、もう少しざっくりとした目標になっており、そこでは全員の目標が決まったと考えていた。しかし、メンバー各々が少しずつ違った目標をイメージしており、今回の目標が決定するまでの間、メンバーの認識に誤差が生まれてしまった。この修正に時間がかかってしまったことも今回のプロジェクトを完成できなかった要因の一つだと考える。

今回の失敗を踏まえ、より良い取り組みを行うために考えたことは2つである。一つ目は、事前準備をもっと行うことである。ざっくりとした目標定義だけでなく、具体的な目標を決め、また、どう言った経緯で各々がどのような動きをするのかをよく話し合う必要があると考えた。二つ目は、メンバー間での話し合いである。動きが決まっているとしたとしても、急な変更や臨機応変に対応していかなければいけない場面はとて多く存在すると、今回の取り組みで実感した。行っている課題の進行状況やアイデアだけではなく、毎回授業初めに目標を確認したりすることで認識の誤差を減らすことが出来るのではないかと考えた。

7 まとめ

今回の実験で学んだことは、情報の共有の重要性である。各々の考えや思いの共有、話し合いを沢山しなかったことで、発生している問題が多々あった。これが出来ていれば、実験目標のズレだけでなく、1人が用事や体調不良で休んだとしても遅れることなく作業が出来、発表でもスムーズに発表することが出来たと考えられる。

参考文献

- [1] YouTube ヘルプ, <<https://support.google.com/youtube/answer/7239739?hl=ja&sjid=14891511027695807438-AP>>, 2023/06/08.
- [2] ForestSeo, Python と YouTubeAPI で急上昇の動画のデータを取ってくる, <<https://qiita.com/g-k/items/7c98efe21257afac70e9>>, 2023/06/08.