

ニコニコデータセットの分析

原田玲央

۲

バイクでのログを用いた分析

北海道行ってきました！！



せっかくだしログでもとれないかな、、、

このアプリ使えそう、、、



シンプル加速度ロガー

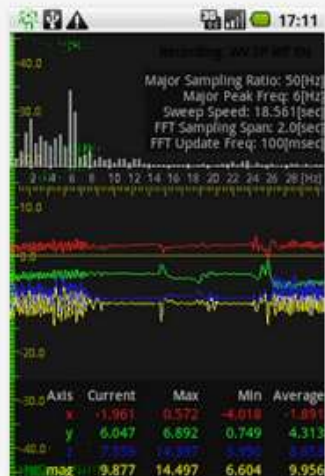
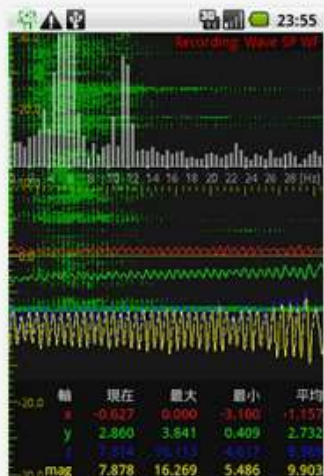
Sora Takayama - 2010年9月9日
ツール

インストール済み

このアプリはお使いの端末に対応しています。

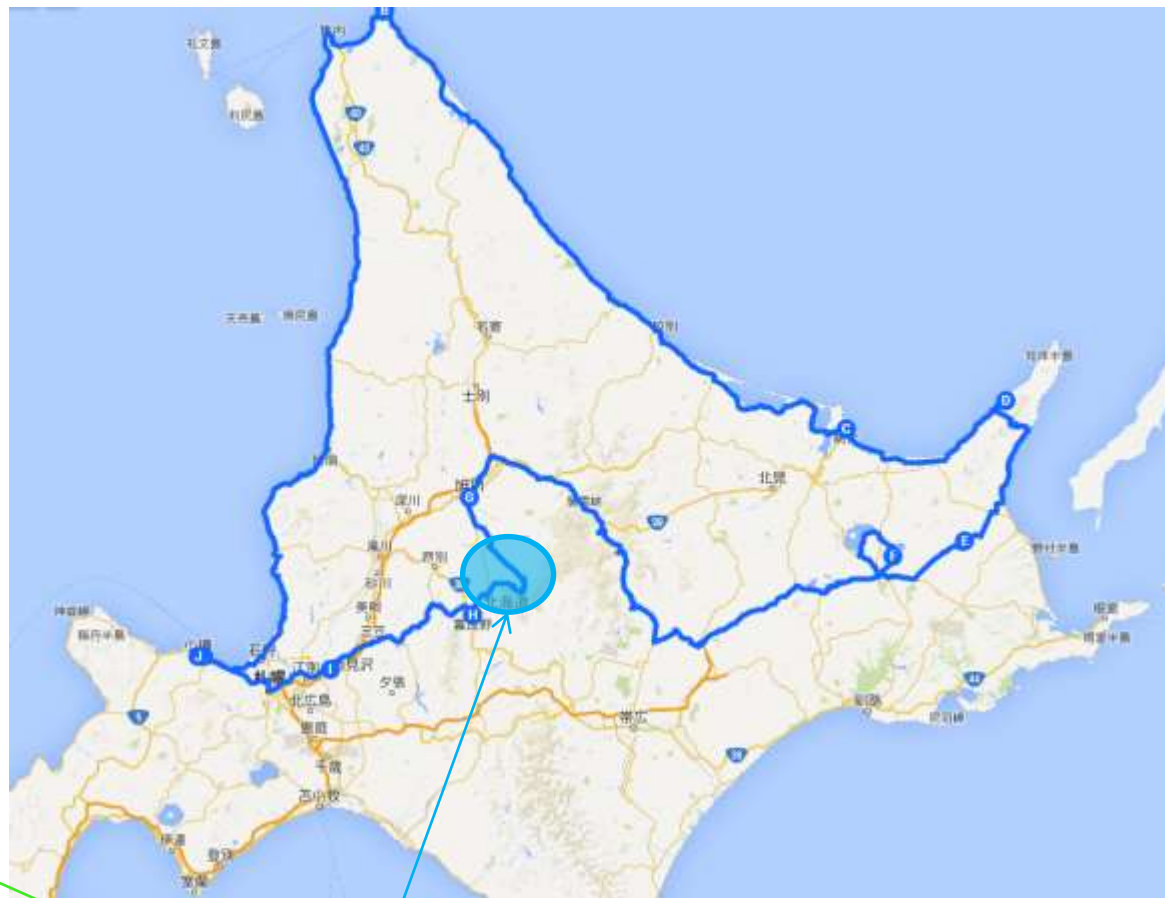
★★★★★ (166)

x,y,z軸に対する加速度を
取得できて、csvファイルと
して出力してくれる





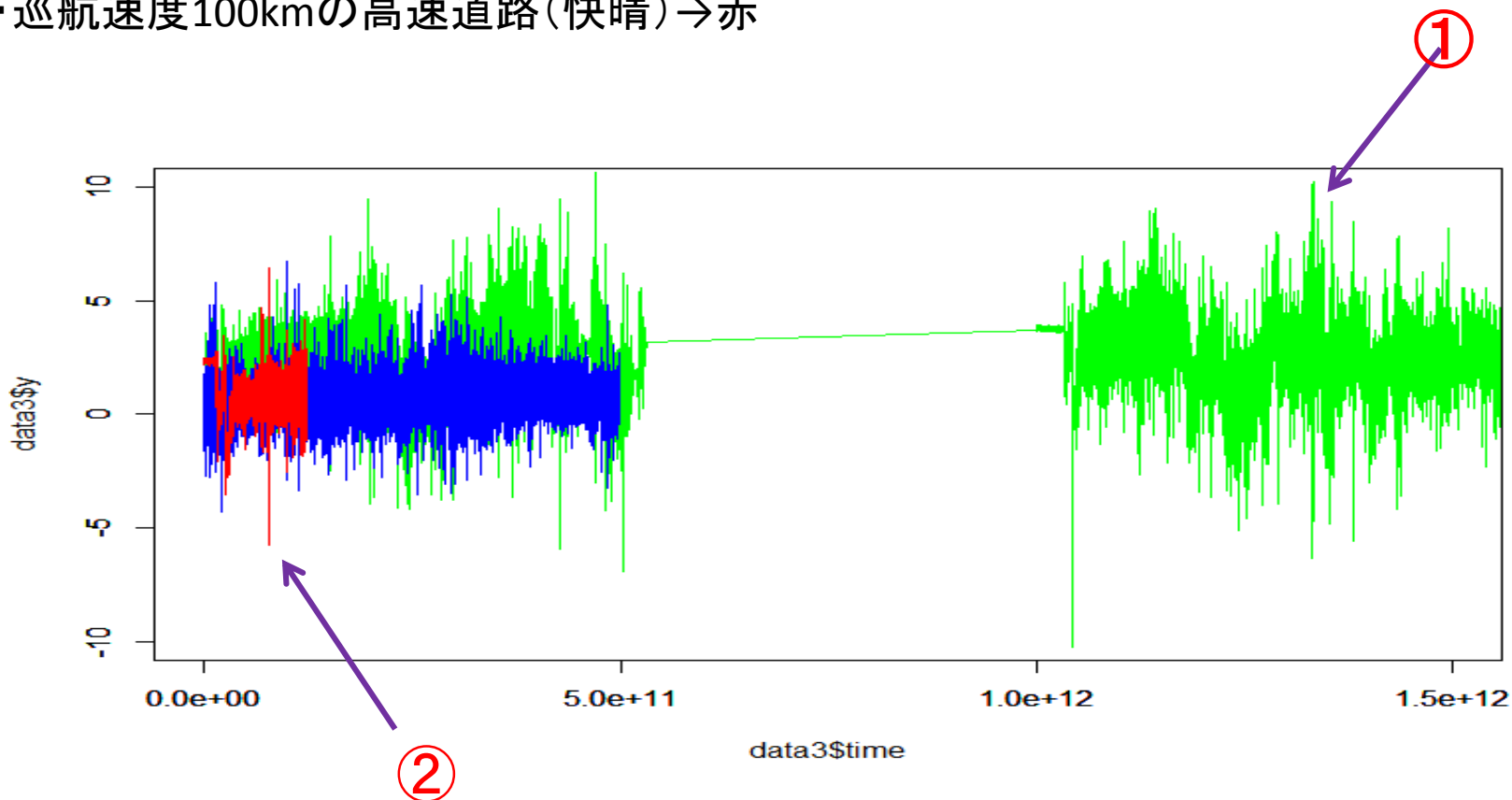
スマートフォンをハンドル部分につけて左右への傾きを計測



3ヶ所でデータを取ってきました！

分析結果

- ・巡航速度60km前後の2車線の峠(快晴)→緑
- ・巡航速度30km～40km前後の一本道の峠(濃霧)→青
- ・巡航速度100kmの高速道路(快晴)→赤



※約15万行のcsv

①巡航速度が比較的速い峠

→左右へのグラフの変動が大きい

②高速道路の合流地点(車線変更)

→加速度がかかる

つまり、

車線変更を頻繁にする車両はこのような波形が多く見られるのでは？

2. ニコニコデータセットを用いた分析

ニコニコデータセットって、、、

- 国立情報学研究所のダウンロードサービスより(株)ドワンゴが提供
- 2007年3月～2012年11月に投稿された約830万件のデータ
- 動画のメタデータ(12GB)
- コメントデータ(300GB)

データの量が膨大すぎる(;´д`)

ということで一部だけ利用しました

とりあえず中身を見してみる



```
{ "video_id": "sm1001", "thread_id": 1179132..... }
  ⇒ { "key": "value", ..... }
```

※key=video_id,thread_id,upload_time,length,size_high,view_counter,tag,etc...

(´・ω・`)。oO(フムフム、、、なるほど、、)

このデータを使って何か分析してみよう

- I .カテゴリータグの出現確率
- II .タグの数と再生数の関係
- III .タグを用いたアソシエーション分析

調べてみた!!

I .カテゴリータグの出現確率

2009年12月18日 01時23分 投稿
再生: 1,181,295 | コメント: 28,749 | マイリスト: 19,739
料理 カテゴリ前日総合順位:181位(過去最高:1位) | この動画の記事
動画の詳細情報を開く

登録済

いいね! 1,703 8+1 5

メニューを開く

登録タグ (16件)

- 料理
- SEGAの人
- ゆっくり解説
- ニコニコ技術部
- ニコニコ挫折部
- 自作PC
- ドリームキャスト
- 発想の勝利
- まっくろくろいの
- ペン・トー
- 黒子ショー

コメント NG設定

通常コメント

コメント	再生時
ラヴィ	01:15 ▲
www	01:15
あらら	02:36
www	02:31
www	02:48

あ...
おいwww
素手www

カテゴリータグ

エンターテインメント、音楽、歌ってみた、演奏してみた、踊ってみた、VOCALOID、ニコニコインディーズ
動物、料理、自然、旅行、車載動画、スポーツ、ニコニコ動画講座、歴史
政治

科学、ニコニコ技術部、ニコニコ手芸部、作ってみた

アニメ、ゲーム、アイドルマスター、東方、ラジオ、描いてみた

例のアレ、日記、その他

R-18

ファッション

まずデータの加工

```
222433 ゲーム ↵
222434 音楽 ↵
222435 ゲーム ↵
222436 ゲーム ↵
222437 ゲーム ↵
222438 踊ってみた ↵
222439 アニメ ↵
222440 ゲーム ↵
222441 料理 ↵
222442 ゲーム ↵
222443 ↵
222444 踊ってみた ↵
222445 ↵
222446 ゲーム ↵
222447 ゲーム ↵
222448 ゲーム ↵
222449 ゲーム ↵
222450 ゲーム ↵
222451 ↵
222452 ゲーム ↵
222453 エンターテイメント ↵
222454 日記 ↵
222455 ゲーム ↵
222456 東方 ↵
222457 その他 ↵
222458 スポーツ ↵
222459 アニメ ↵
222460 演奏してみた ↵
222461 歌ってみた ↵
222462 歌ってみた ↵
222463 ゲーム ↵
222464 ゲーム ↵
222465 動物 ↵
222466 ゲーム ↵
222467 音楽 ↵
222468 ゲーム ↵
222469 例のアレ ↵
222470 VOCALOID ↵
222471 音楽 ↵
222472 音楽 ↵
222473 音楽 ↵
222474 ゲーム ↵
222475 ゲーム ↵
```

Json形式のファイル

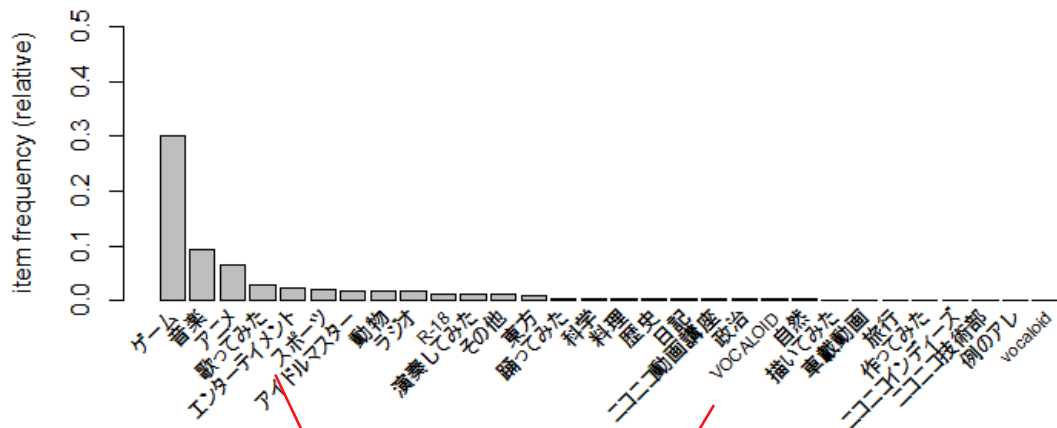


カテゴリタグだけを抽出



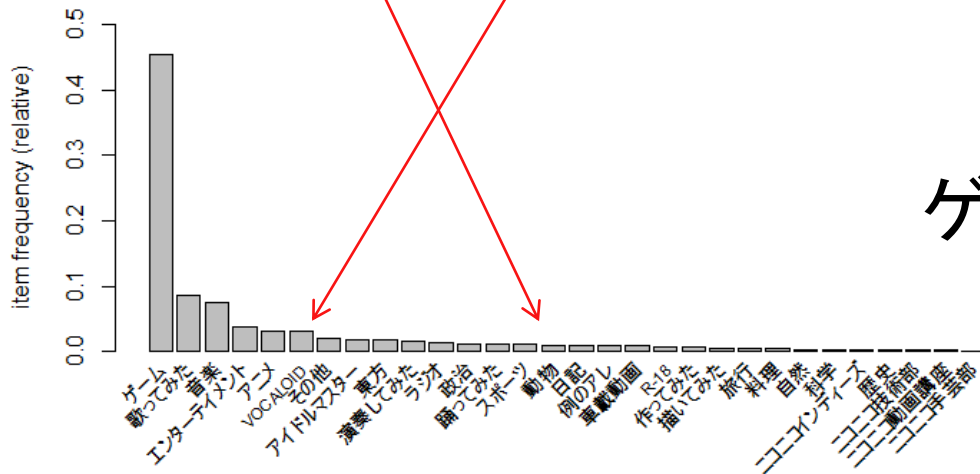
RでitemFrequencyPlot(data)

～2008/01



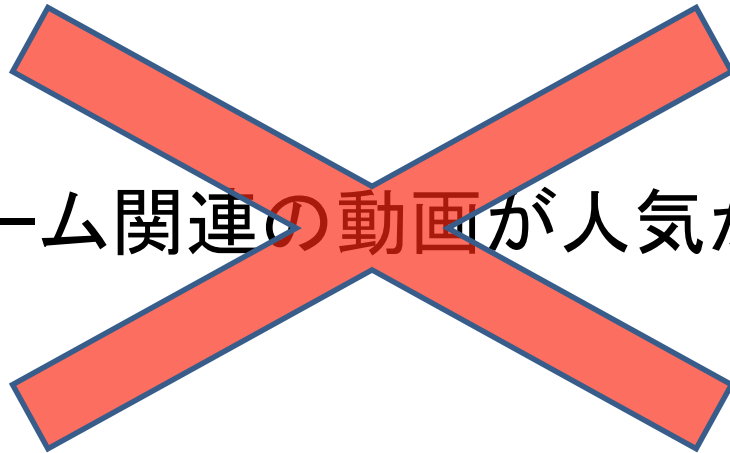
- ・ゲームタグが多く出現
- ・VOCALOIDタグが増加
- ・スポーツタグが減少

2012/06～2012/11

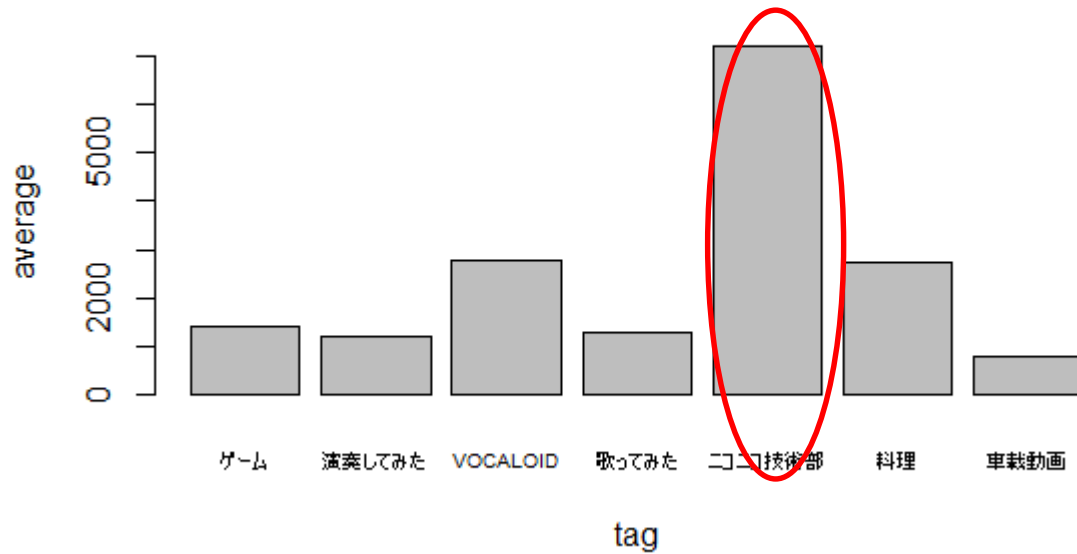
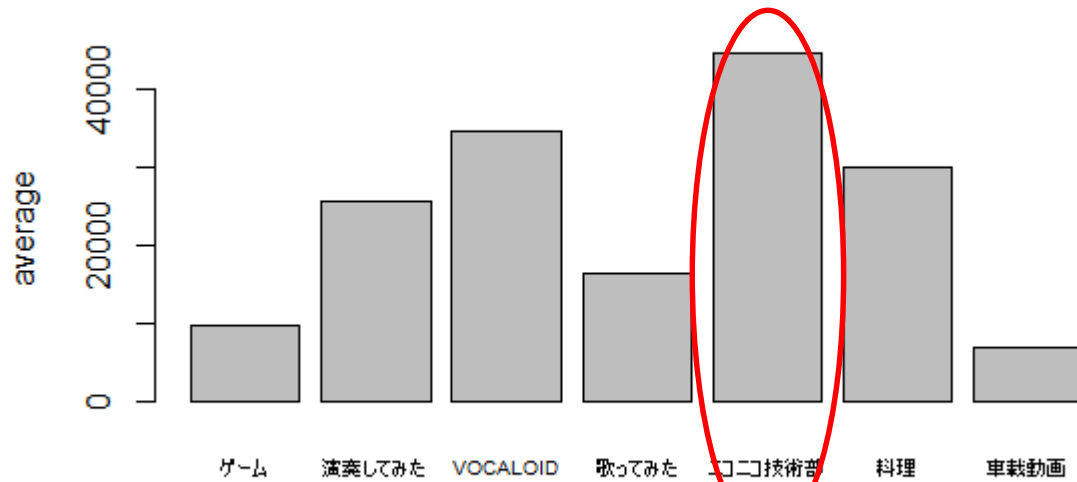


ゲームタグが3割超！！

(´・ω・`).oO(ゲーム関連の動画が人気があるのか、.)



ある期間におけるタグごとの平均再生回数



この2つの期間においては共に「ニコニコ技術部」が人気？

投稿する側

→ゲーム関連の動画をよくup

視聴者側

→(※この期間においては)「ニコニコ技術部」の動画をよく再生
(=おもしろい動画が多い)

Ⅱ. タグの数と再生数の関係

- タグは1つの動画に対して最大11個
(うちカテゴリタグは3個まで)
- キーワードを含むタグ検索が可能

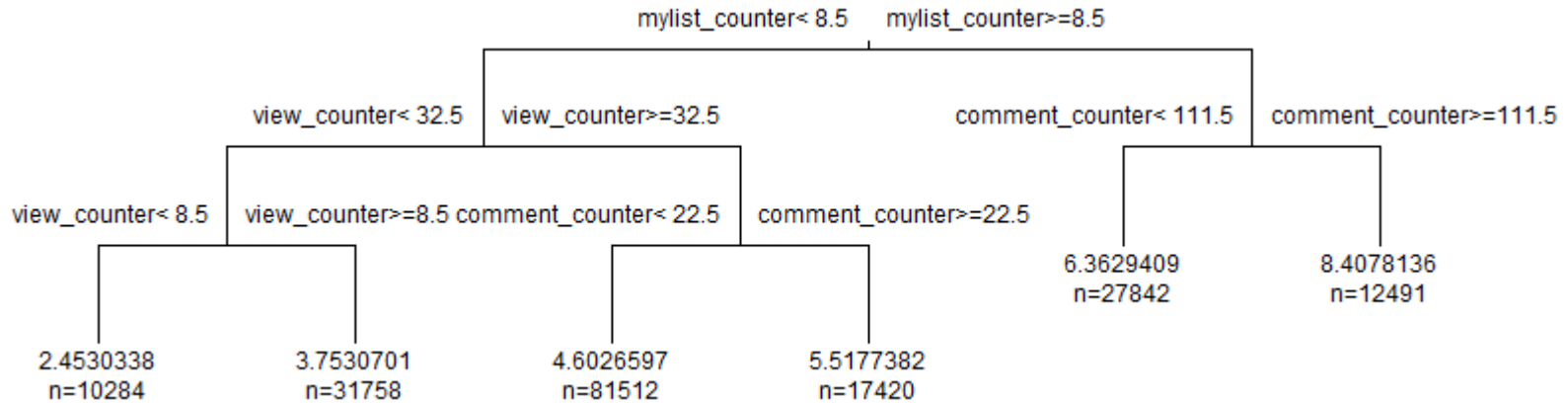


タグが多いほど盛り上がっている動画なのでは？

まずデータの加工

	A	B	C	D	E	F	G	H	I	J	K	L
1	video_id	thread_id	upload_time	title	movie_type	size_low	size_high	mylist_counter	length	view_counter	comment_counter	tag_counter
2	sm19247769	1351682087	2012-10-31T20:14:47+09:00	夕暮れ時に聴	mp4	0	334397195	4	6370	285	3	7
3	sm19190025	1351008005	2012-10-24T01:00:04+09:00	ゴールデンボ	mp4	0	330772456	105	6068	4794	819	9
4	sm19056468	1349519736	2012-10-06T19:35:36+09:00	【ラテール】【1	mp4	0	329617200	6	6052	392	2	4
5	sm19009272	1348966550	2012-09-30T09:55:50+09:00	HITS! THE T	mp4	0	314243418	6	5758	959	3	5
6	sm19136152	1350398270	2012-10-16T23:37:50+09:00	ゴールデンボ	mp4	0	307575580	100	5643	4119	304	10
7	sm19222191	1351398493	2012-10-28T13:28:13+09:00	10/28朝礼拝	mp4	0	301622420	0	5528	3	0	6
8	sm19030839	1349186606	2012-10-02T23:03:26+09:00	ゴールデンボ	mp4	0	295198172	109	5416	11931	333	11
9	sm19180689	1350912067	2012-10-22T22:21:07+09:00	【2005年】永	mp4	0	287891470	13	5278	1284	49	1
10	sm19224809	1351417903	2012-10-28T18:51:43+09:00	10/28夕礼拝	mp4	0	285215448	0	5227	7	0	6
11	sm19082678	1349791487	2012-10-09T23:04:47+09:00	GAREI ZERO	mp4	0	282847214	3	5856	82	55	4
12	sm19083775	1349796714	2012-10-10T00:31:54+09:00	ゴールデンボ	mp4	0	280324231	127	5145	9968	467	11
13	sm19023945	1349099321	2012-10-01T22:48:41+09:00	悪魔のQP	mp4	0	280178882	65	5155	11829	87	10
14	sm19167701	1350774072	2012-10-21T08:01:12+09:00	HITS! THE T	mp4	0	279514768	3	5122	627	5	6
15	sm19055425	1349511144	2012-10-06T17:12:24+09:00	MADOKA MA	mp4	0	278997793	37	5164	574	13	4
16	sm19171937	1350819360	2012-10-21T20:36:00+09:00	10/21夕礼拝	mp4	0	278162582	0	5099	3	0	6
17	sm19135257	1350393578	2012-10-16T22:19:38+09:00	映画史上最悪	mp4	0	272720968	22	4998	774	71	10
18	sm19121900	1350224647	2012-10-14T23:24:07+09:00	【走行音】ホ	mp4	0	272057736	2	5025	85	3	5
19	sm19169045	1350792945	2012-10-21T13:15:45+09:00	10/21朝礼拝	mp4	0	263064774	0	4822	2	0	5
20	sm19238090	1351568249	2012-10-30T12:37:29+09:00	うちの時間	mp4	0	262028267	0	8060	18	0	3
21	sm19127648	1350303905	2012-10-15T21:25:05+09:00	ふいあ通	mp4	0	260418835	6	5341	199	2	5
22	sm19098172	1349997416	2012-10-12T08:16:56+09:00	【うんこちゃん	mp4	0	252414169	11	4665	2493	22	1
23	sm19039298	1349288965	2012-10-04T03:29:25+09:00	LAST EXILE	mp4	0	249814624	14	4669	359	2	3
24	sm19237958	1351565832	2012-10-30T11:57:12+09:00	うちの時間	mp4	0	246462591	0	7583	1	0	4
25	sm19115627	1350183938	2012-10-14T12:05:38+09:00	HITS! THE T	mp4	0	242030175	1	4435	329	4	6
26	sm19225207	1351421737	2012-10-28T19:55:37+09:00	スフィアのオ	mp4	0	238194567	18	6120	463	20	2
27	sm19018052	1349019648	2012-10-01T00:40:47+09:00	切ない癒し	mp4	0	237465340	22	4525	633	6	3

- Jsonからメタデータ抽出＋タグの数をカウント
- csvとして出力
- rpartで決定木(目的変数＝tag_counter)

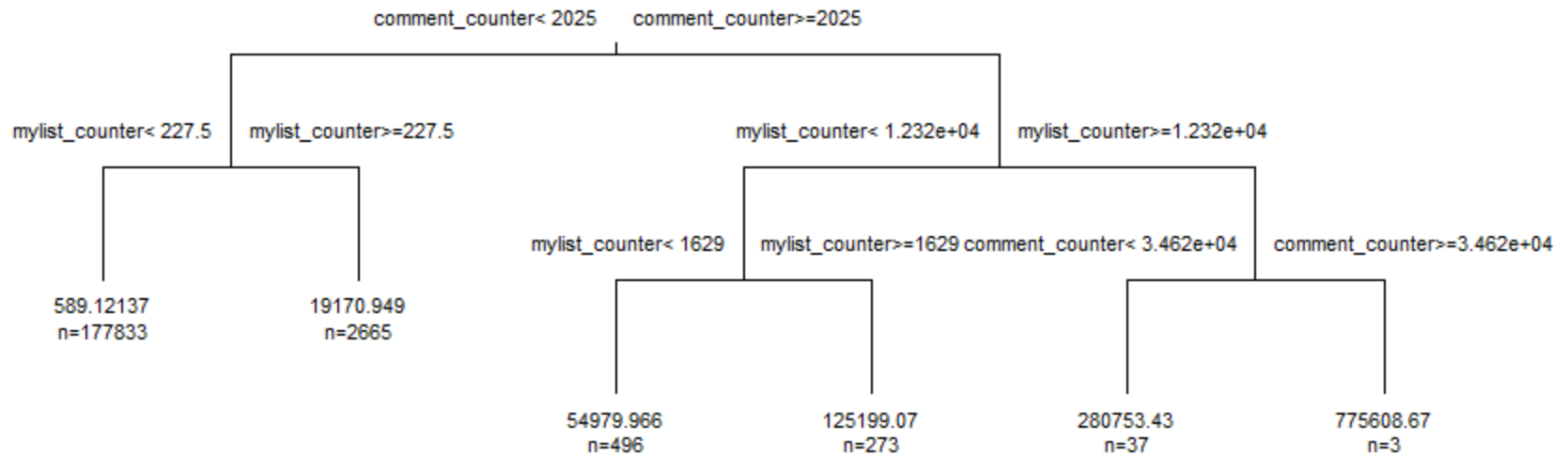


- ・分岐基準はマイリス登録数、再生回数、コメント数
- ・マイリス登録が多くコメントも多い→タグの数も多い
- ・マイリス登録が少なく、再生数も少ない→タグの数も少ない

しかし、
分岐基準となる値が小さい→少しでも見られている動画はタグの数が少ないことはない？

つまり、

盛り上がる動画はタグの数が多くつけられる
逆は??



目的変数 = view_counter

分岐基準にtag_counterが出てこない



再生数が多いのはタグの数によるとは限らない

Ⅲ. タグを用いたアソシエーション分析

タグの間にどのような関連性があるのかを調べてみた

条件

2012年8月、9月に投稿されたデータ(35万件)

その期間のうち「ニコニコ技術部」タグを含む(818件)

まずデータを加工

```
775 user/10691693,くだらないニート動画,ニコニコ手芸部,ニコニコ技術部,またお前か,作って,紳士の社交場,紳士向け+
776 UTAU,UTAUプラグイン配布有り,UTAU支援ツール,その他,ニコニコ技術部,もっと評価されるべき+
777 1/1690タクタンク,F-34,ウオッカ,コミッサール,ソビエト連邦軍歌,ソ連軍主力戦車,ニコニコ技術部,傑作戦車ロージナ,赤いニコニコ動画,赤外線バトルタンク,
778 3008,lightnave,ニコニコ技術部,パンク風いてない(確認),ゆっくり,ゆっくり解説,一時停止非推奨,触手+
779 エンターテイメント,くだらないニート動画,ニコニコ技術部,作って,人生精んでる,作って,兵器,武器+
780 エンターテイメント,くだらないニート動画,どうしてこうなった,ニコニコ技術部,人生精んでる,作って,兵器,武器+
781 カワサキパロウのグループ,Shit'D'it'aka(か)!!!,ただし魔法は民から出る,ダイソウの申し子なのに,ニコニコ技術部,ニコニコ格闘同好会,仮面ライダー
782 555,シーク注意,タイマー10,タイマー10 555,ニコニコ技術部,電子工作+
783 sm17961005,sm18139422,sm18877815,sm19242253,どうしてこうなった,ニコニコ技術部,作って,作って,実験,科学+
784 くだらないニート動画,どうしてこうなった,ニコニコ技術部,やっみた,人生精んでる,作って,科学+
785 くだらないニート動画,ニコニコ技術部,作って,人生精んでる,作って,実験,科学,転載+
786 Nice Train,じこはおこるさ,デジビュー,ニコニコ技術部,ブラレール,改造ブラレール+
787 LEGO,タモリ倶楽部,ニコニコ技術部,もっと評価されるべき,レゴ+
788 カイジ,どうしてこうなった,ニコニコ技術部,作って,作って,科学+
789 ニコニコ技術部,ニコニコ造形部,フィギュア,フルスクラッチ,ゆっくり解説,作って,東方,東方クラフト,東風早苗,絶対許早苗,鎌田吾作+
790 32000,ニコニコ技術部,俺の愛車,自動車+
791 PC,オーバークロック,でっという,ニコニコ技術部,作って+
792 ※都合のいいタグは全てa主がつけてます,※都合の悪いタグは全てa主がつけてます,「自称」ニコニコ技術部,どうしてこうなった,ドライアイス,ニコニコ
793 サムネ見てからクリック余裕でした,ドイツ,ニコニコ技術部,ベルリンの壁,ロシア,自動ジャンプ,自動車,転載厨,約リ+
794 WOODLOD10オルゴールアレンジ,オルガンニート,オルゴール,カゲロウデバイス,ニコニコ技術部,作って+
795 オシロスコープ,ニコニコ技術部,先行+sm18949683,東方,読技術,転載+
796 L P C 1 7 6 B , L P C X p r e s s o , た こ ん 力 , た こ ん 力 技 術 部 , ニ コ ニ コ 技 術 部 , 電 子 工 作 +
797 ニコニコ技術部,ニコニコ技術部養成講座,薄電接審判,電子工作+
798 Nalvo・Amazon,クラッチオペレーティングシステム,それとなく,ニコニコ技術部,まさかの高1,レストラン,器用首芝,整備,車,車載動画+
799 LED,しくみ,ニコニコ技術部,ニコニコ技術部養成講座,ひかひか,回路,投稿者コメント,解説して,電子工作+
800 どう見てもアンパンマン,ニコニコ技術部+
801 Higa,M/S少女,サラダバー!!!,ニコニコ技術部,プラモデル,まどマギ,ライザーさやか+
802 BSM,Perfume,てってて~,ニコニコ技術部,フリスク,レーザー,レーザービーム(Perfume),作って,実験,電子工作+
803 くだらないニート動画,ニコニコ技術部,作って,人生精んでる,作って,実験,気任,科学+
804 どうしてこうなった,ニコニコ技術部,作って,人生精んでる,作って,実験,科学+
805 user/10691693,※都合のいいタグは全てa主がつけてます,どうしてこうなった,ニコニコ手芸部,ニコニコ技術部,作って,作って,工
806 ※都合のいいタグは全てa主がつけてます,どうしてこうなった,ニコニコ手芸部,ニコニコ技術部,作って,人生話んでる,作って,天才とバカの境界
807 ニコニコ技術部,ヘッドフォン推奨,疑似サラウンド,科学,音楽+
808 iPad,ShareRockPerc,WOODLOD,アプリ,ニコニコ技術部,初音ミク,音楽ゲーム+
809 user/10691693,どうしてこうなった,ニコニコ手芸部,ニコニコ技術部,人生精んでる,卵,科学+
810 user/10691693,エンターテイメント,オナホ,どうしてこうなった,ニコニコ技術部,作って,人生精んでる,作って+
811 BMS,BDF2012,schranz,コントローラー,ニコニコ技術部,作って,専コン,真実次郎,自作コントローラ+
812 アバター,アバター作成ツール,さらばん.com,ニコニコ技術部,マウス総,投稿者コメント+
813 LEGO,カートリッジ式レゴム銃,ゴム銃,ニコニコ技術部,もっと評価されるべき,レゴ+
814 R(ライフル) -18,カンスミス,スリングショット,ニコニコ兵器開発局,ニコニコ技術部,ライフル型スリングショット,作って,職人の仕事,言い値で買お
```

Json形式のファイル



条件に適したタグを抽出
(1行につき1件の
トランザクションデータ)



>read.transactions(data)
>apriori
>inspect

分析結果

```
>library(arules)
```

```
>ap<-apriori(data,parameter=list(maxlen=4,support=0.02,confidence=0.7,ext=T))
```

```
>inspect(head(sort(ap,by="lift"),n=50))
```

```
1 {やってみた, 科学, 作ってみた} => {どうしてこうなった}
```

support	confidence	lhs.support	lift
0.02	0.89	0.02	17.020807834

```
2 {やってみた, 科学} => {どうしてこうなった}
```

0.02	0.85	0.025	16.305647841
------	------	-------	--------------

```
3 {ニコニコ技術部, やってみた, 科学} => {どうしてこうなった}
```

0.02	0.85	0.025	16.305647841
------	------	-------	--------------

```
4 {ニコニコ動画講座} => {ニコニコ技術部養成講座}
```

0.03	0.81	0.039	16.210365854
------	------	-------	--------------

```
5 {やってみた, 作ってみた} => {どうしてこうなった}
```

0.037	0.83	0.045	15.938403520
-------	------	-------	--------------

やってみた、作ってみた⇒どうしてこうなった
というルールが多く出現

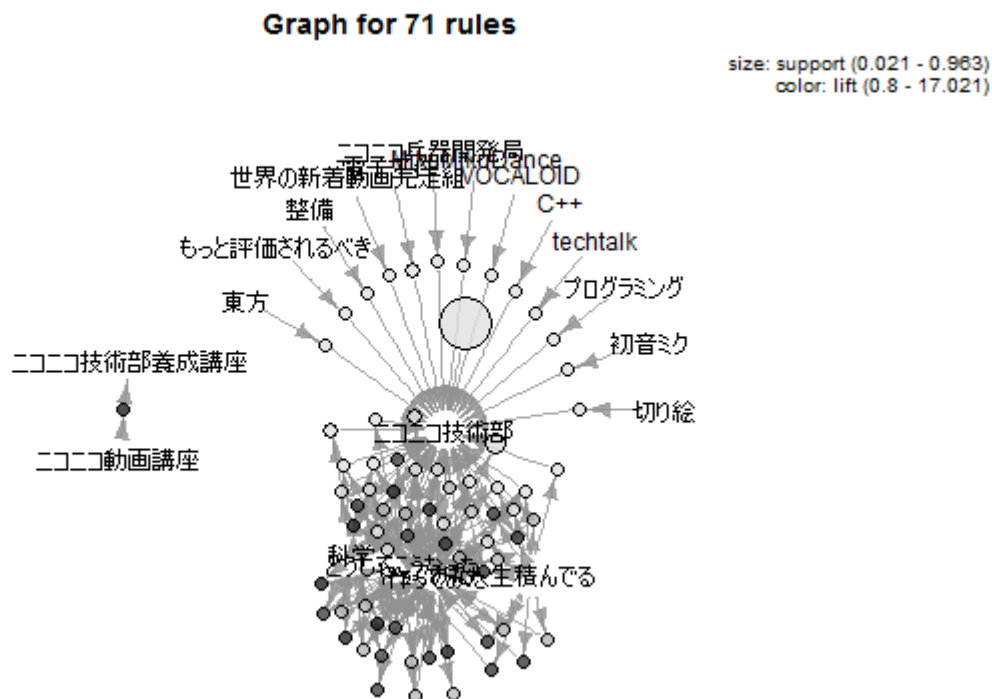


「すごい」というより「面白い」という要素が強そう

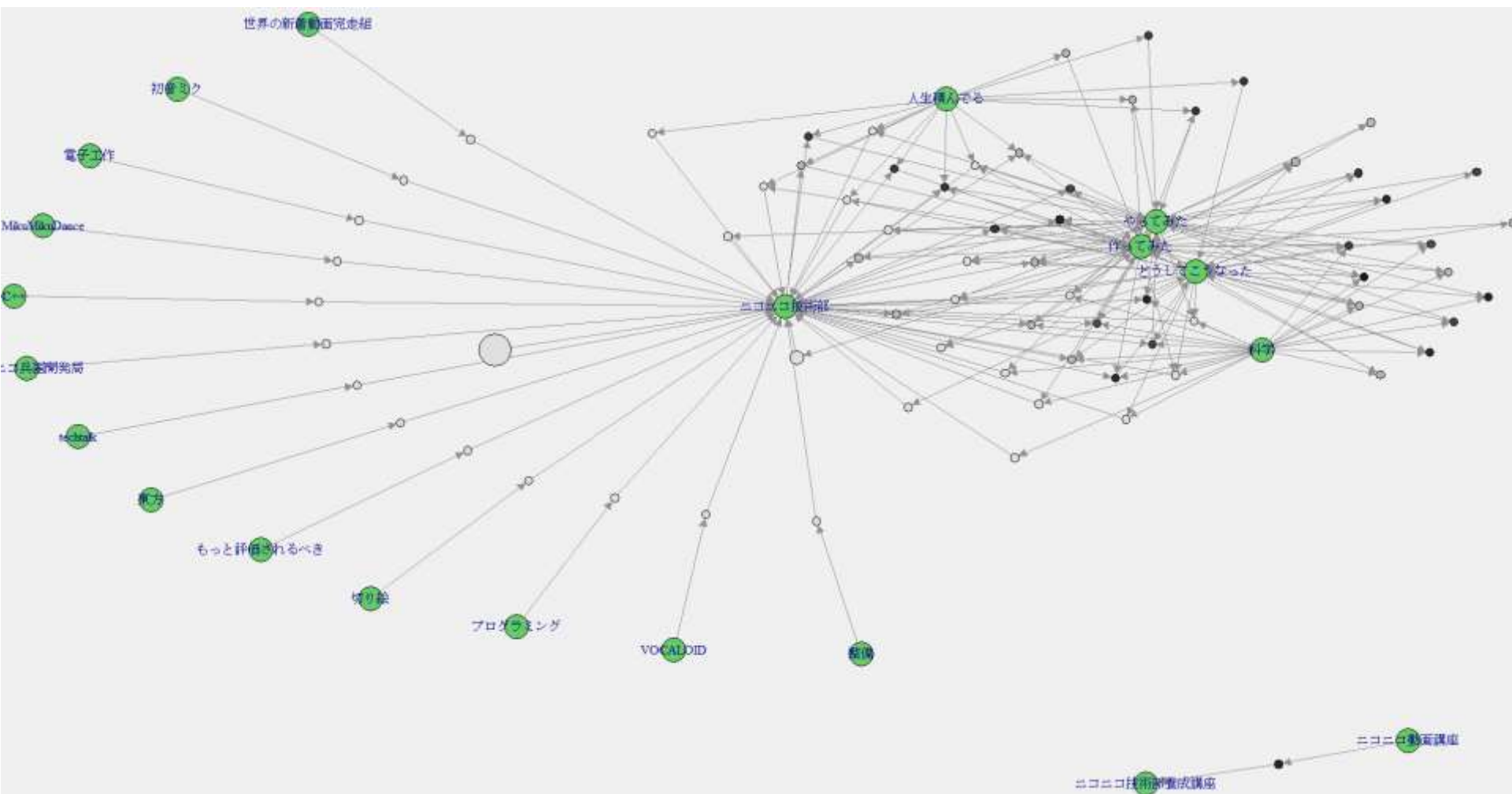
可視化すると、、、

```
>library(igraph)
```

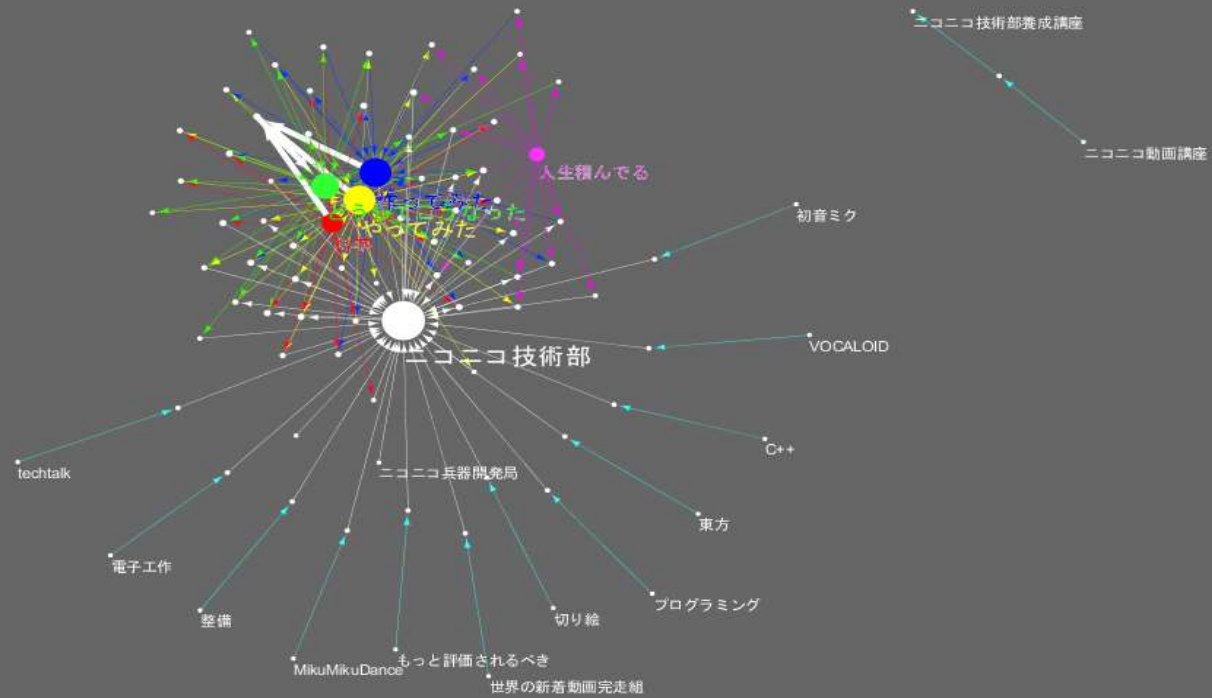
```
>plot(ap,method='graph',control=list(type='items',layout=layout.fruchterman.reingold))
```



Plotにinteractive=TRUEを追記すると ちょっと見やすくなる



Cytoscapeで可視化



わからないということが、わかった

反省点

- コメントデータを使えなかった

→関連性の薄いタグが複数あったら「タグ理解」のコメントが出現するのでは？

→wwwが複数あったら8888がある可能性が高いのでは？

ご清聴ありがとうございました