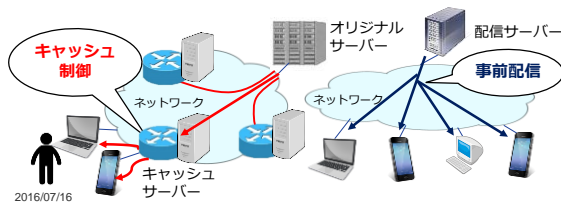


ユーザー生成コンテンツの 視聴数推移パターン分析と 人気推移予測

田中 達也† 阿多 信吾‡ 村田 正幸†
† 大阪大学 大学院情報科学研究科
‡ 大阪市立大学 大学院工学研究科

研究背景

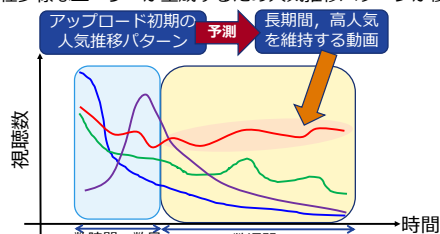
- YouTube に代表されるユーザー生成コンテンツ (UGC ; User Generated Content) の視聴の普及
- **人気コンテンツを早期かつ高精度に判別**することが有効
 - ネットワークトラフィック削減のための適切な**キャッシュ制御**
 - 配信サーバのピーク時負荷抑制のための**事前配信**
 - 効果的な**広告マーケティング**



2016/07/16

研究課題

- UGCの将来の人気度予測は**困難**
- 多種多様なユーザーが生成するため人気推移パターンが複雑



アップロード初期の人気推移パターンから
長期間、高人気を維持する動画を予測できないか？

2016/07/16

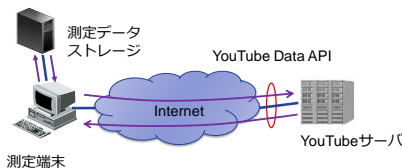
研究の目的と方法

- 研究目的
 - UGC の人気推移パターンの分析
 - **アップロード初期の時点**で将来高人気を維持する動画の予測
- 研究の手順
 - YouTube 動画の視聴数時系列データの収集
 - 人気度推移の分析
 - クラスタリングによる視聴数推移パターン分析
 - ▶ **k-means 法**を用いたクラスタリング
 - 人気度予測と評価
 - ▶ 教師あり学習の一種である**単純ベイズ分類器**により判別

2016/07/16

YouTube データの概要

- 新着動画の視聴数の時系列データ
 - YouTube Data API version3 [13] を用いて取得した動画
 - ▶ (動画数 : 87,830 2015/10/14~2015/12/16)
 - アップロード 1 週間までの 1 時間毎の視聴数
 - アップロード 1 週間経過後の 1 日毎の視聴数

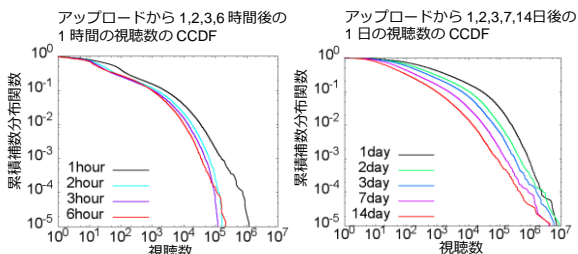


[13] "YouTube Data API" <https://developer.google.com/youtube/v3/>

2016/07/16

アップロードから経過時間による視聴数分布

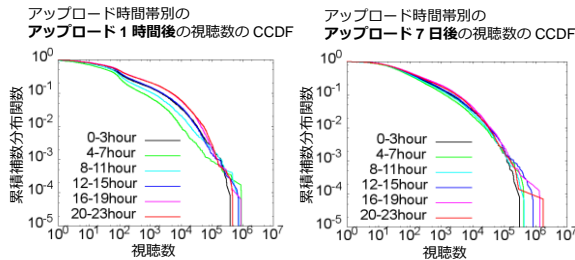
- アップロード直後が最も視聴数が多く、その後減少
- 両対数グラフで裾野の部分が線形に近いカーブ
 - 非常に視聴数の高い少数の動画が存在



2016/07/16

アップロード時刻による視聴傾向の差異

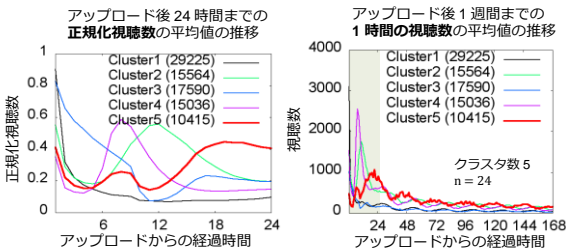
- 世界標準時刻で16-19時、20-23時の視聴数が多い
- **アップロード直後**はインターネット人口の多い時間帯にアップロードされた動画が多く視聴数を獲得する傾向



2016/07/16

クラスタリングによる視聴数推移パターン分析

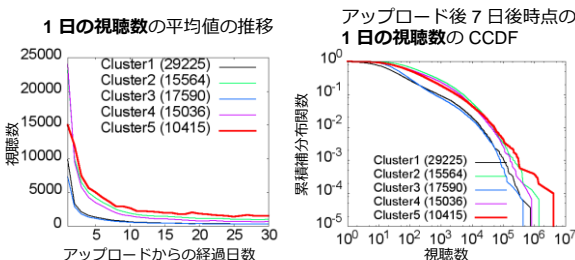
- k-means 法を用いて初期の視聴数推移パターンを分類
- 各動画に対して、最初の n 時間の視聴数の最大値で各時間の視聴数を割った正規化視聴数をもとに k-means 法を適用
- **視聴数を維持する推移パターン (クラスタ 5) の存在を確認**



2016/07/16

各クラスタの人気度推移

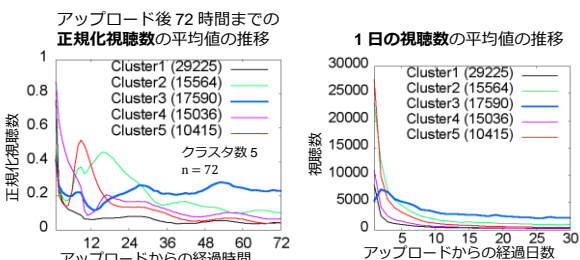
- 初期の時点で正規化視聴数が維持されているクラスタは、その後も他のクラスタより視聴数が高い傾向
- **クラスタ 5 が視聴数を維持**



2016/07/16

クラスタリングに用いる期間の変更

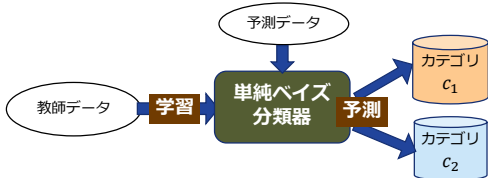
- n=48,72 で k-means 法を用いた場合
- n=24 の場合と同様の傾向を確認
- n=72 のときクラスタ 3 が **視聴数を維持する推移パターン**



2016/07/16

単純ベイズ分類器 (NBC : Naive Bayes Classifier)

- ベイズの定理を用いた**教師あり学習**による分類法
- **学習**: 教師データを用いて各カテゴリ c の学習サンプルがある入力セット f_1, \dots, f_n を有する確率を計算
- **予測**: f_1, \dots, f_n を有する各予測サンプルを事後確率が最大になる c に分類
 - ▷ $\text{classify}(f_1, \dots, f_n) = \arg \max_c p(C=c) \prod_{i=1}^n p(F_i = f_i | C=c)$



2016/07/16

単純ベイズ分類器の学習と予測の方法

- 入力
 - アップロード初期 Y 時間の視聴数の最大値で各時間の視聴数を割った**正規化視聴数**と**視聴数の最大値の桁数**
- 出力
 - C1 : **高人気維持** C2 : それ以外
 - ▷ C1 の定義 1 : d 日後までの d 日間の累積視聴数が全体の上位 1%
 - ▷ C1 の定義 2 : d 日後の 1 日の視聴数が全体の上位 1%

学習データの例

動画 ID	Y 内の正規化視聴数				Y 内の最大視聴数の桁数	分類カテゴリ
	スロット1	スロット2	...	スロットY		
abcdefghijkl	1.0	0.5	...	0.4	5	C1
lmnopqrstuv	0.5	0.2	...	0.0	3	C2
wxyz1234567	0.2	0.8	...	0.6	4	C1

2016/07/16

