

PD3中間報告会

4EP4-3S 鶴瀬和輝

研究テーマ
YouTube動画の再生数を予測する

背景

- ・YouTubeは、2019年広告売上高150億ドル、月間アクセス20億人に達するほど、世界で最も普及した動画マーケティングツール
- ・再生数を増やすための施策が、日々多くのメディアでノウハウが紹介されている(動画広告、セミナー、記事)

背景と課題

再生数を増やすための施策が、日々多くのメディアでノウハウが紹介されている

各取り組みによりどれだけ視聴回数を伸ばすことが出来るのか、定量的データに基づいたノウハウ紹介は少ない、何か効果的のかがついていない

ゴール

YouTube(動画クリエイター)にとって生命線とされる再生数を伸ばすために、どのような施策をいけるかを知る必要があるのかを知り、動画の再生数を予測できるモデルを構築する。

アプローチ・方法

使用するデータ：Prospace YouTube動画視聴回数予測
コンペティションのデータ
レコード数：19720
列数：17

変数名	説明
video_id	動画ID(一意のID)
title	動画タイトル
duration	動画の長さ(秒)
channel	チャンネルID
upload_date	アップロード日時
category	動画のカテゴリ
keywords	動画のキーワード
likes	動画のいいねの数
views	動画の再生回数
dislikes	動画のいいねと反対の数
comment_count	動画のコメントの数
comment_dislike	動画のコメントのいいねと反対の数
rating	動画の平均評価
tags	動画のタグ

アプローチ・方法

EDA(探索型データ解析)を行い、再生数を予測するのに効果的な特徴量を見つける。

- 以下の特徴量を用意
- ・dislikes(低評価数)、likes(高評価数)comment_count(コメント数の四捨五入)log_2(二乗などのaggregation(集計)特徴量
 - ・duration、likes、comment_count(ラベル)
 - ・テキスト：(channel_title, category)などこのaggregation特徴量
 - ・テキストのtf-idf → svd、doc2vec、tf-idf+svd
 - ・Target Encoding
 - ・テキストの中にkeywordあり、なしのbinary特徴量

アプローチ・方法

- ・モデリングでは、LightGBMという勾配ブースティングフレームワークを用いる

なぜLightGBMを使うのか？

- ・カテゴリ変数に対して特別な処理を自動的に実行してくれるので、One-Hotエンコーディングの手間を無くせる
- ・既存のデータセットを極力加工せずに利用するという観点で、特徴量エンジニアリングの負担を軽減してくれる

進捗状況

- 特徴量について
- ・データ収集、Prospace YouTube動画視聴回数予測コンペ
 - ・収集したデータをEDA(探索型データ解析)を行う
 - ・特徴量エンジニアリング
 - ・numerical data(dislikes, likes, comment_count)のaggregation特徴量を作成
 - ・text dataはtf-idf特徴量をsvd特徴量に変換、そしてaggregation特徴量を作成
 - ・Target Encoding
- モデルについて
- ×・シングルモデルでは、LightGBM、XGBoost、CatBoostを試す
 - ×・各々のモデルのパラメータ調整

今後の進め方について

- ・BERTを使ってテキストデータ(title, description, channel_title, tags)の特徴抽出
- ・過剰状態でできていなかったモデリングを行う