

PD3中間報告会

4EP4-35 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、セクシー、起訴)

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増やすことが出来るのか、正確な予測を導き出すことが出来れば、広告効果は上がる、何が効果的かわかっていない

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、どのような施策を打てば再生数が増えるのかを知り、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

変数名	変数の説明
video_id	動画ID(主キー)
title	動画タイトル
description	動画説明
channel_title	チャンネル名
category	チャンネルカテゴリー
upload_date	動画アップロード日時
likes	動画のいいねの数
dislikes	動画のいいねの数
comment_count	動画のコメントの数
view_count	動画の再生回数
tags	動画のタグ
is_live	動画がライブか否か
is_private	動画が公開されているか
is_unlisted	動画が公開されているか
is_short	動画がショート動画か否か
is_monetized	動画が収益化されているか
is_family_safe	動画が家族向けか否か
is_ad	動画が広告か否か

使用するデータ:
Produce YouTube動画視聴回数予測コンペ
コンペディションのデータ

レコード数: 19720
列数: 17

アプローチ・方法

EDA(探索型データ分析)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。

- 以下の特徴量を使用
- video_id(主キー)、likes(いいねの数)、comment_count(コメントの数)、view_count(再生回数)
 - channel_title、category、upload_date、tags、is_live、is_private、is_unlisted、is_short、is_monetized、is_family_safe、is_ad
 - テキストのtag → svd、docvec、tfidf+one
 - 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?

- カテゴリー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減らせる
- 既存のデータセットを効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、Produce YouTube動画視聴回数予測コンペ
 - 収集したデータをEDA(探索型データ分析)を行う
 - 特徴量エンジニアリング
 - numerical data(likes, dislikes, comment_count)のaggregation特徴量を作成
 - text dataはtag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
 - × 各々のモデルのハイパーパラメータ調整

今後の進め方について

- ・ BERTを使ってテキストデータ(title, description, channelTitle, tags)の特徴抽出
- ・ 書籍状態で使えていなかったモデリングを行う