https://www.kaggle.com/c/kkbox-music-recommendation-challenge/data?select=members.csv.7z
If you are running archlinux: • git clone https://aur.archlinux.org/kaggle-api.git • cd kaggle-api • makepkg -si • Go to the first link and create a kaggle account and agree to the competition rules • go to your account page on kaggle and create an api key and save the kaggle.json file in the folder ~/.kaggle/ • kaggle competitions download -c kkbox-music-recommendation-challenge Import required dependencies
import numpy as np import pandas as pd import tensorflow as tf import tensorflow.data as tfd from loguru import logger from tqdm import tqdm from typing import List, Any, Tuple, Optional, Dict [191 datapath: str = os.path.join('', 'data')
def list_files(directory: str, extension: str) -> List[str]: all files = os.listdir(directory) return [os.path.join(directory, file) for file in all_files if file.split('.')[-1] == extension] datasets: Dict[str, pd.DataFrame] = dict() for filepath in tqdm(list_files(datapath, 'csv'), ascii=True, desc="Loading data from disk"): datasets[os.path.basename(filepath).split('.')[0]] = pd.read_csv(filepath) Loading data from disk: 100% ###################################
Some key metrics on dataset size [194. print(f"Loaded {len(datasets)} csv files") for key, value in datasets.items(): print(f"length of dataset {key} is {len(value)}") Loaded 6 csv files length of dataset songs is 2296320 length of dataset test is 2556790 length of dataset sample_submission is 2556790 length of dataset sample_submission is 2556790 length of dataset train is 7377418
length of dataset members is 34483 length of dataset song_extra_info is 2295971 [195 # Let's remove the 'sample_submission' dataset _ = datasets.pop('sample_submission') Below is a more detailed analysis of the raw data, omitted for brevity [1]: # for key, value in datasets.items(): # print(f"Information for dataset: {key}") # print(f"Description") # print(f"Description") # print(f"Description")
print(\n') # print(\n') # print(\n') # print(\dataframe 'head'") # print(\dataframe 'head'") # print(\n') # print(\n') # print(\n') # print(\n') "English' Description of the Dataset • we have a list of users, their personal information, the songs that they liked and didnt like, and where they accesed the song • we also have metadata about each song in the dataset Dataset descriptions:
 members. csv user ID, location, gender, song_extra_info + songs.csv song_id, song name, artist name, composer, lyricist, genre id, song length, language train.csv: userid, songid, where they found the song, whether they liked it or now This will be used to create clusters of songs, or song groups
Recommendation dataset: • userid, song rating, songs cluster, songs feature vector Below is a preview of the raw curated 'songs' dataset [198- raw = pd.DataFrame()
raw["genre_id"] = datasets["songs"]["genre_ids"] raw["song_length"] = datasets["songs"]["song_length"] raw["song_length"] = datasets["songs"]["song_length"] raw["song_id"] = datasets["songs"]["song_length"] raw["song_id"] = datasets["songs"]["song_length"] raw [199
2 SUPER JUNIOR NaN V465 31.0 231781 DWVVVurlrpuz+XPUpFvucctVQEyPqcpUkHR0ne1RQzPs0= 3 S.H.E 湯小陳 徐世珍 465 3.0 273554 dKMBWozyScdxSkihKG-YV47nc18N9q4m58+b4e7dSSE= 4 貴族精選 Traditional Traditional 766 52.0 140329 W3bqWd3T+VeHFZHAUTARgW9AvVRaF4N5Yzm4Mr6Eo/cs=
2296319 Kim Robertson Nan Nan 829 52.0 270466 V/9qPSUOGT0FelGBhDlZRR25Y5yyjdXC2Z/+bJExo5s= 2296320 rows × 7 columns Here we begin cleaning the curated songs dataset for analysis • This includes: • removing na values • removing unicode characters
<pre></pre>
raw['artist_name'] = raw[raw['artist_name'].apply(only_letters)]['artist_name'] raw['lyricist'] = raw[raw['lyricist'].apply(only_letters)]['yricist'] raw['composer'] = raw[raw['omposer'].apply(only_numbers)]['composer'] raw['genre_id'] = raw[raw['genre_id'].apply(only_numbers)]['genre_id'] We are working with a massive dataset, and to allocate enough memory to run an MCA on even just raw cat would take over 1.4 TiB so lets drop everything with an NA [203. raw.dropna(axis=0, how='any', inplace=True)
Taw = raw[(raw.alanguage != -1] One other way to reduce the number of categories for MCA is to filter our dataframe by the top 1000 artists [286 num_artists=10e
raw = raw[raw.artist_name.isin(raw['artist_name'].value_counts(sort=True).index.tolist()[:num_artists])] raw = raw[raw.composer'].value_counts(sort=True).index.tolist()[:num_posers])] raw = raw[raw.lyricist:], value_counts(sort=True).index.tolist()[:num_posers])] raw = raw[raw.genre_id.isin(raw['genre_id'].value_counts(sort=True).index.tolist()[:num_posers])] raw = raw[raw.language.isin(raw['language'].value_counts(sort=True).index.tolist()[:num_posers])] Below is a look at the artists, composers, lyricists, languages, and genres being preliminarily analyzed [210. raw.artist_name.unique()
array(['erica', 'Prince', 'Fromm', 'MISIA', 'Vitas', 'millstones', 'Mili',
'Rio', 'Nucc', 'Djavan', 'Tellement', 'Shura', 'narry', 'Ken', 'BOKKA', 'nyanyannya', 'Inmay', 'Gotttod', 'PiA', 'DELPHIC', 'SKALAPPER', 'Sarantos', 'Megapuss'], dtype=object) [211. array(['nao', 'Prince', 'Fromm', 'Sinkiroh', 'Vitas', 'millstones', 'Mili', 'keeno', 'YUI', 'MUNA', 'Jain', 'Frost', 'Chumbawmba', 'toku', 'ryo', 'Scandal', 'Cloa', 'erica', 'UZ', 'YUKA', 'UVERworld', 'Oohyo', 'Myk', 'Owlle', 'mothy', 'Mot', 'TAKUYA', 'Cranfield', 'wintermute', 'Hizaki', 'Kayoko', 'Soko', 'mito', 'MYX', 'MINATO', 'Gil', 'Ostesun', 'ZEEBRA', 'Takuro', 'Mami', 'Yk', 'MINATO', 'Gil', 'Ostesun', 'ZEEBRA', 'Takuro', 'Mami', 'Yk',
'cocco', 'ALEX', 'Toraboruta', 'minato', 'Mogwai', 'TAKURO', 'Royz', 'Vanomai', 'VaLSHE', 'Cozzi', 'Miwa', 'PENTCILLIN', 'Chara', 'Shuho', 'Lisa', 'kous', 'Austra', 'E', 'Takuya', 'DIV', 'MASTERWORKS', 'Babyface', 'INORAN', 'Eumlansonyeon', 'Rina', 'SYK', 'RUKA', 'Miliyah', 'Jiro', 'Traditional', 'Hisashi', 'SYK', 'RUKA', 'Miliyah', 'Jiro', 'Traditional', 'Hisashi', 'HAKUEI', 'Fandacchi', 'Sakoshin', 'Halstead', 'Windycat', 'Cocco', 'yellowhato', 'tofubeats', 'dustbox', 'BackHoEdic', 'PollyphonicBranch', 'Neru', 'Charm', 'Chay', 'Marmozets', 'AnDisM', 'KAMIJO', 'Depapepe', 'Valo', 'Tyuryu', 'Estra', 'Kamijo', 'KAMIJO', 'Depapepe', 'Valo', 'Tyuryu', 'Estra', 'Kamijo', 'Shinkiroh', 'Mastodon', 'YoOKEY', 'MOGWAI', 'Misia', 'Heavenz', 'Aquibird', 'Soil', 'Wowaka', 'Garry', 'Chay', 'Candyman', 'Tatsuhi', 'SHIBU', 'Baekja', 'Nhato', 'Korn', 'INOVADER', 'FACT', 'Rika', 'DATEKEN', 'Teru', 'Feldberg', 'RYOSPEED', 'Kreativgang', 'RIka', 'DATEKEN', 'Teru', 'Feldberg', 'RYOSPEED', 'Kreativgang', 'LMPSH', 'TERU', 'Geoon', 'Haruna', 'OHLO', 'HISASHI', 'OOHYO',
'Rio', 'Tomomit', 'YUKIYOSHI', 'Kozhevnikov', 'GISHO', 'CHISA', 'MASASHI', 'Jin', 'ORYO', 'YUKKE', 'Djavan', 'Taishi', 'HIZAKI', 'Tellement', 'doriko', 'KNOXX', 'Isaac', 'Shura', 'narry', 'Chiquewa', 'Ken', 'Bokka', 'NoRIKIYO', 'dezzy', 'GODES', 'Dmitrii', 'nyanyannya', 'Inmay', 'Issac', 'Tetoteto', 'MAH', 'Rado', 'Chiquewa', 'Miya', 'Gotttod', 'heroism', 'Scarecrow', 'Pla', 'Yk', 'Delphic', 'myremy', 'Tatsuro', 'LeE', 'Kember', 'HIDE', 'Yuya', 'NAGMATIC', 'nodoka', 'chobi', 'Cypher', 'toraboruta', 'Versailles', 'SHINCO', 'T', 'Shirai', 'Yohan', 'JIGG', 'AnDism', 'satoshi', 'Hajimetal', 'Taizo', 'Sarantos', 'Atsushi', 'DolToKi', 'mabanua', 'KOYAMNUSIC', 'SATOchi', 'Megapuss', 'PATRON', 'Makira', 'Kentaro', 'TeddyLoid', 'Simz', 'Shinji', 'Joi', 'KeN', 'mooumoon', 'Saltwater', 'shogo', 'Clankidd', 'Rhymerube', 'uz'], dtype=object)
raw.lyricist.unique() array(['erica', 'Prince', 'Fromm', 'Misia', 'Vitas', 'millstones', 'Mili',
'Chiquewa', 'Halstead', 'Windycat', 'Tablo', 'Cocco', 'chelly', 'Mami', 'toflubeats', 'dustbox', 'SALU', 'PolyphonicBranch', 'Neru', 'Charm', 'CHISA', 'Chay', 'Marmozets', 'AnDisM', 'KAMIJO', 'Depapepe,' 'Valo', 'Yuzvu', 'Mastodon', 'MoGMAI', 'Heavenz', 'Aquibird', 'SOil', 'Wowaka', 'chay', 'Candyman', 'Ryuji', 'SHIBU', 'Baekja', 'Mhato', 'Korn', 'toku', 'FACI', 'Pika', 'DATEKEN', 'Feldberg', 'RYOSPEED', 'Kreativgang', 'LMPSH', 'JIRO', 'TERU', 'eaeon', 'Haruna', 'HISASHI', 'OOHYO', 'Rio', 'minato', 'ORYO', 'Tatsuro', 'Djavan', 'Shogo', 'Taishi', 'VICTiM', 'Tellement', 'RUCCA', 'Isaac', 'MINATO', 'Shura', 'narry', 'chobi', 'Keen', 'BOKka', 'dezzy', 'GODEES', 'Nayanannya', 'Imay', 'Isac', 'Chiquewa', 'GODEES', 'Nayanannya', 'Imay', 'Isac', 'Miya', 'LEE', 'Keember', 'Teru', 'HIDE', 'Nodoka', 'Kyo', 'Caboruta', 'ShiRai', 'flansowa', 'Vohan', 'AnDism', 'Haj metal', 'Sarantos', 'Megapuss, 'Pandacchi', 'Kentaro', 'TeddyLoid',
Scaling the song_length to the range [0, 1]. This is done because without this, the FAMD only discriminates based off of song lenth and tends to ignore the other factors in the analysis [215. from sklearn.preprocessing import MinMaxScaler [216. scaler = MinMaxScaler() raw["song_length"] = scaler.fit_transform(raw["song_length"].values.reshape(-1, 1)) raw [217. from sklearn.preprocessing import MinMaxScaler [218. artist_name composer lyricist genre_id language song_length song_id
Record R
228954 PIA
An MCA (Multiple Correspondence Analysis) is performed on all categorical variables A PCA (Principal Component Analysis) is performed on all continuous variables A PCA is performed on the outputs of the previous MCA and PCA, and the results of this are returned We pruned the dataset so heavily becuase MCA requires us to one-hot-encode our categorical variables. Becuase we are dealing with hundreds of classes and thousands of samples, the memory overhead is tremendous. The memory overhead of the entire dataset is 1.4TB! from prince.famd import FAMD famd = FAMD(
check_input=True, engine='auto', random_state=42) famd = famd.fit(raw.drop(['song_id'], axis='columns')) [218.
118 2.191764 -0.435968 0.502185 1125 0.18833 -0.748530 0.505635 1742 -0.989763 1.241269 -0.434505 1852 1.861413 -0.472608 0.570857 2288152 -0.91648 0.165201 -0.629373 2289054 0.083141 -0.373721 -0.539500 2293048 1.780624 -0.909563 0.588362
2294539 -0.935966 0.164797 -0.536861 2816 rows × 3 columns Below is a plot showing the results of the FAMD • From this we can see that you need more than just song information to make a prediction, as a large number of songs are mostly indistinguishable based off of their metadata alone [219- ax = famd.plot_row_coordinates(
<pre>X = raw, ax=None, figsize=(20, 20), x_component=0, y_component=2, color_labels=["Artist: {}".format(t) for t in raw["artist_name"]], ellips=_outline=False, ellips=_fill=True, show_points=True } Row principal coordinates • Artist: Aquibird</pre>
Artist: Babyface Artist: Babyface Artist: Babyface Artist: Candyman Artist: Chiquew Artist: Chiquew Artist: Chiquewanha Artist: Chumbawanha Artist: Cocco Artist: Cocco Artist: Corafield Artist: DEPAIEN
Artist: EdolST Artist: EdolST Artist: Elsmlansonyeon Artist: FaCT Artist: Feldberg Artist: Fromm Artist: Fromm Artist: From Artist: Fromt Artist: GARNIOELIA Artist: GARNIOELIA Artist: GHANIOELIA Artist: Hill Artist: Hill Artist: Heavenz Artist: Heavenz
Arist: INORAN Arist: Immay Arist: Immay Arist: Immay Arist: Immay Arist: Immay Arist: Immay Arist: Ken Arist: Ken Arist: Ken Arist: Ken Arist: Ken Arist: Ken Arist: Kes Arist: Lead Arist: Lead Arist: Lead Arist: Lead Arist: Liad Arist: Liad Arist: Liad Arist: Liad Arist: Liad Arist: Liad Arist: Mannanta
Artist: Magauss Artist: Millyah Artist: Milyah Artist: NiGritMaRe Artist: NiGritMaRe Artist: NiGritMaRe Artist: Oorlivo Artist: Oorlivo Artist: Oorlivo Artist: Oille Artist: PilA
Artist: Polyphonic Branch Artist: Roy Artist: Roy Artist: Sov Artist: SAINPAPER Artist: SCAINDAL Artist: SCAINDAL Artist: SCAINDAL Artist: SCAINDAL Artist: SFAIAR Artist: SFAIAR Artist: SFYAIR Artist: SFYAIR Artist: SFYAIR Artist: Sarantos Artist: Siondive Artist: Slowdive Artist: Slowdive Artist: Slowdive
-2 -1 0 1 2
Artist: trica Artist: millstones Artist: milva Artist: mothy Artist: mouron Artist: narry Artist: ryanyannya Artist: tyanyannya Artist: tyanyannya Artist: tyuny Artist: tyunya Artist: super-ell Artist: tofubeats Artist: wintermute
ax.get_figure().savefig('famd_row_coordinates_100_artists_song_length.svg') Gradient Boosting Our classifier created using the information discovered above Imports [234
Read in the data [275 sei = pd.read_csv('/data/song_extra_info.csv') members = pd.read_csv('/data/members.csv', parse_dates=['registration_init_time', 'expiration_date']) songs = pd.read_csv('/data/songs.csv') train = pd.read_csv('/data/train.csv') test = pd.read_csv('/data/test.csv')
members['expiration_year'] = members['expiration_tatte'].dt.year members arron (columns = ['registration_init_time'], 'expiration_date'], inplace = True) members['members['expiration_year'] = members['registration_date'].subtract (members['registration_date'].subtract (members['registration_date']).dt.days.astype(int) members['registration_year'] = members['registration_vear'] = members['registration_init_time'].dt.year members.head() [257.
3 mCuD+tZthERA/oSGPqk38e041J8ZSBaLcu7nGolivhle 1 0 NaN 9 1 2015 2015 4 q4HRBIVSssAFS9iRtxWohxuk9kCYMKjHOEagUMV6rQ= 1 0 NaN 4 138 2017 2017 Create a valid training and test set from combining the avalilable dataframes [258- train = train.merge(songs , on='song_id' , how='left') train = train.merge(sei , on = 'msno' , how='left') train = train.merge(sei , on = 'song_id' , how='left') test = test.merge(songs , on='song_id' , how='left') test = test.merge(members , on = 'msno' , how = 'left') test = test.merge(members , on = 'msno' , how = 'left')
test = test.merge(sei , on = 'song_id' , how = 'left') 258. 3116 259. test['song_length'].fillna(test['song_length'].mean() , inplace = True) test['song_length'] = test['song_length'].astype(np.uint32) test['language'].fillna(test['language'].mode().values[0] , inplace= True) test['language'] = test['language'].astype(np.int8) train['song_length'].fillna(train['song_length'].mean() , inplace = True) train['song_length'] = train['song_length'].astype(np.uint32) train['language'].fillna(train['language').mode().values[0] , inplace= True)
<pre>train['language'] = train['language'].astype(np.int8) print(train.columns) train.head() False</pre>
'composer', 'lyricaist', 'languager', 'city', 'bd', 'gender', 'registration_year', 'name', 'isrc', 'year'], 'dtype='object') msn song_id source_steen_nam source_treen_nam source_treen_name source_treen_nam source_treen_na
3 Xumu+NIjS6QYVxDS4/t3SawvJ7viT9hPKXmf0RtLNx8= 2A87tzfnJTSWqD7gIZHisolhe4DMdzkbd6LzO1KHjNs= my library Local playlist more local-playlist 1 255512 1019 Soundway Kwadwo Donkoh 13 24 female 9 2301 2011 2017 Disco Africa GBUQH1000063 2010. 4 FGtllVqz18RPiwJj/edr2gV78zirAiY/9SmYvia+kCg= 3qm6XTZ6MOCU11x8FIVbAGH5l5uMkT3/ZalWG1oo2Gc= explore Explore playlist 1 187802 1011 Brett Young Brett Young Relly Archer Justin Ebach 1 0 NaN 7 2103 2012 2017 Sleep Without You QM3E21606003 2016.
train['genre_ids'].fillna('no_genre_id' , inplace= True) train['genre_ids_count'] = train['genre_ids'].apply(lambda x: genre_count(x)).astype(np.int8) test['genre_ids'].fillna('no_genre_id' , inplace= True) test['genre_ids_count'] = test['genre_ids'].apply(lambda x: genre_count(x)).astype(np.int8)
train['genre_ids'].fillng('no_genre_id', implace= True) train['genre_ids'].fillng('no_genre_id', implace= True) train['genre_ids'].fillng('no_genre_id', implace= True) test['genre_ids'].fillng('no_genre_id', implace= True) test['genre_ids'].fillng('no_genre_id', implace= True) test['genre_ids'].fillng('no_genre_id', implace= True) test['genre_ids'].fillng('no_genre_ids').apply(lambda x: genre_count(x)).astype(np.int8) false
Train['genre_ids'].fillna('no_genre_id', inplace= True)
property of the state of the st
Part
The content of the
Part
Part
Part
Company Comp
The content of the
Continue
The content of the
18
Companies Comp
The content of the
The content of the
The content of the