



Data Analytics

What's your accent ?



Table of Contents

INTRODUCTION	3
CONTEXT	3
OBJECTIVES.....	3
BUSINESS CASES	4
PROJECT TOOLS	6
WORKFLOW.....	7
DATA & DATA SOURCES	8
KAGGLE.....	8
SPEECH ACCENT ARCHIVE	10
DATA COLLECTION.....	10
KAGGLE (<i>REST API</i>)	10
SPEECH ACCENT ARCHIVE (<i>Web Scraping</i>).....	13
DATA CLEANING & EXPLORATORY DATA ANALYSIS	13
DATA CLEANING	13
EXPLORATORY DATA ANALYSIS	14
X.....	Error! Bookmark not defined.
DATABASE TYPE SELECTION	14
RELATIONAL DATABASES.....	14
NON-RELATIONAL DATABASES.....	14
DATABASE CREATION	16
SQL QUERIES.....	16
Cloud: BigQuery.....	17
ENTITY RELATIONSHIP DIAGRAM	18
GENERAL DATA PROTECTION REGULATION (GDPR)	19
TYPE OF DATA.....	19
API DEVELOPMENT.....	19
CREATING REST API WITH FLASK.....	19
FOLDERS STRUCTURE	19
YAML FILE	20
APP1.PY	21
WEBSITE	22

HTML	22
CSS	22
MACHINE LEARNING	23
MODELS SELECTION	23
MODELS APPLICATION	23

INTRODUCTION

CONTEXT

Whether you are facing customers in the Tourism Industry, the Luxury field, the restaurants management, or any other front office activity, understanding your interlocutor's culture is paramount. It is of essence, not just for a more efficient communication, but also to better serve and avoid misinterpretation.

OBJECTIVES

The goal of this project is to investigate the relationships between native language and English accent. A step beyond is to predict native language based on the accent.

BUSINESS CASES

The ability to predict a speaker's native language based on their accent in English could have several potential business applications, for example in:

- **Language Learning Platforms:** Companies offering language learning platforms could utilize this technology to provide personalized learning experiences for users. By analyzing their accent in English, the platform could identify their native language and tailor the learning materials, exercises, and pronunciation practice specifically for that language background.
- **Customer Service and Support** Businesses with international customer bases could use this technology in their customer service and support departments. By identifying the native language of the caller based on their accent, companies can route calls to agents who speak the caller's native language, improving communication and customer satisfaction.
- **Market Research and Targeted Advertising:** Companies conducting market research or targeted advertising campaigns could use this technology to gather insights about the demographics of their audience. By analyzing the accents of respondents or customers, businesses can infer their native languages and adjust their marketing strategies accordingly to better resonate with specific language groups.
- **Language Assessment and Testing:** Educational institutions or organizations offering language proficiency assessments and testing could

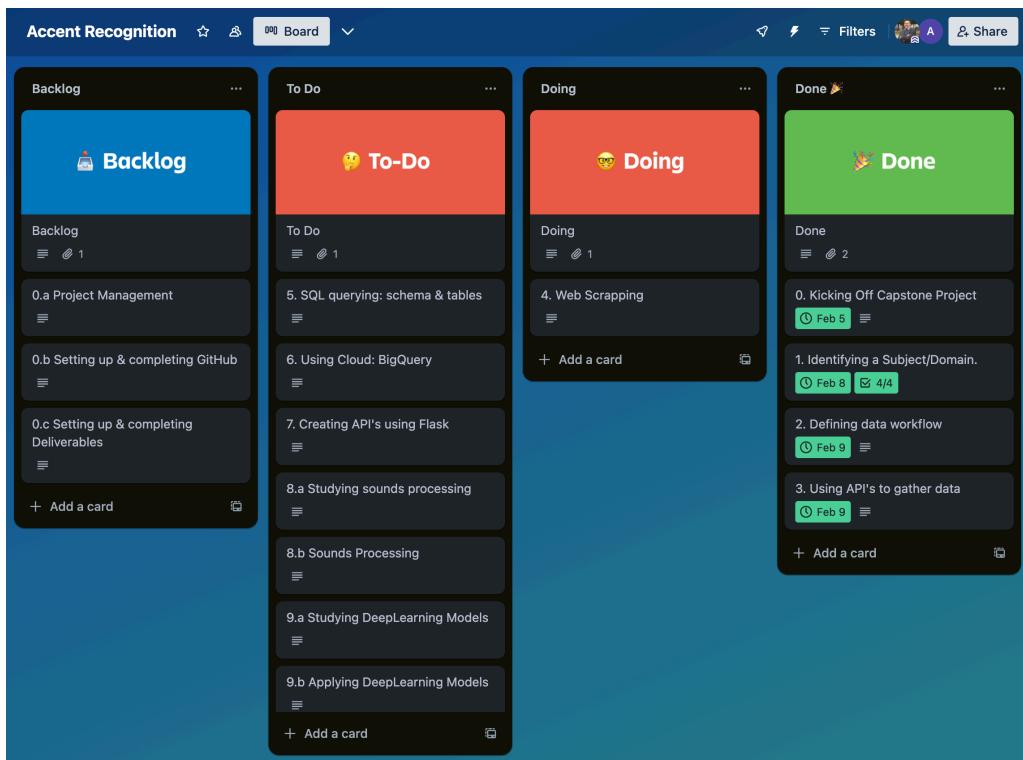
incorporate this technology to improve the accuracy of their assessments. By analyzing the accents of test takers in English, institutions can better understand their language backgrounds and adjust scoring or evaluation criteria accordingly.

- **Speech Recognition and Natural Language Processing:** Companies developing speech recognition and natural language processing technologies could benefit from this capability by enhancing the accuracy and performance of their systems. Understanding the native language of the speaker can help improve speech recognition accuracy and enable more effective natural language understanding, especially in multilingual environments.
- **6. Cross-Cultural Communication and Training:** Organizations with diverse international teams or conducting cross-cultural business could use this technology for communication and training purposes. By understanding the native languages of team members based on their accents, companies can provide targeted cross-cultural communication training and resources to improve collaboration and understanding among team members.
- **7. Immigration and Border Control:** Government agencies responsible for immigration and border control could leverage this technology for various purposes, such as identifying the native languages of individuals during immigration interviews or border crossings. This information could assist in providing language-specific assistance or processing procedures.

- **Content Localization and Translation:** Companies involved in content localization and translation services could use this technology to streamline their processes. By identifying the native languages of content creators or clients based on their accents, businesses can better allocate resources for translation and localization efforts, ensuring high-quality and culturally relevant content for target audiences.

PROJECT TOOLS

1. Trello (Kanban board)



2. Git & GitHub

WORKFLOW

1. Retrieving Data from **Kaggle & Speech Accent Archive**:
 - a. Using Kaggle's API
 - b. Web scraping Speech Accent website
2. Loading datasets into BigQuery for initial data wrangling
3. SQL querying, to create a schema and normalized tables in MySQL
4. Loading data to MySQL
5. Loading data into python dataframe
6. Conducting Exploratory data analysis & Visualization
7. Building a REST API with flask, exposing data from MySQL database
8. Machine Learning:
 - a. Selecting Models
 - b. Performing feature selection & feature engineering
 - c. Applying models

DATA & DATA SOURCES

For this project we are using 2 data sources: Kaggle and Speech Accent Archive
KAGGLE

The data exported from Kaggle consists of a CSV, a text file, and mp3 recordings:

- speakers_all.csv
- reading-passage.txt
- recordings (folder with mp3's)

The data in “reading-passage.txt” is depicted below:

Field Name	Description
age	Speaker's age at the time of the recording.
age_onset	Age at which the speaker started learning english.
birthplace	Location the speaker was born.
filename	Name of the mp3 file (speaker's recording).
country	Country where the speaker was born.
file_missing?	Indicates whether the speaker's recording is available or not

Empty column	Empty column
Empty column	Empty column

The data in “reading-passage.txt” contains the text that each speaker reads outloud:

Description
Please call Stella. Ask her to bring these things with her from the store: Six spoons of fresh snow peas, five thick slabs of blue cheese, and maybe a snack for her brother Bob. We also need a small plastic snake and a big toy frog for the kids. She can scoop these things into three red bags, and we will go meet her Wednesday at the train station.

The “recordings” folder hosts **2134** mp3 files

Description
afrikaans1.mp3 afrikaans2.mp3

SPEECH ACCENT ARCHIVE

Data collecting through webscraping

DATA COLLECTION

KAGGLE (*REST API*)

Data gathered from **Kaggle** has been captured through **Kaggle's API & Kaggle's CLI**.

Kaggle documented their API through the yaml file KaggleSwagger.yaml that can be found here: <https://github.com/Kaggle/kaggle-api/blob/main/KaggleSwagger.yaml>

Hence, in conjunction with **Swagger UI editor**, it helps explore and navigate through Kaggle's API structure.

The screenshot shows a GitHub repository interface. On the left, there's a sidebar titled "Files" with a dropdown menu set to "main". Below it is a search bar and a "Go to file" input field. A list of files is displayed under the "kaggle" directory, including ".gitignore", "CHANGELOG.md", "CONTRIBUTING.md", "KaggleSwagger.yaml" (which is highlighted), "KaggleSwaggerConfig.json", "LICENSE", "MANIFEST.in", "README.md", "setup.cfg", and "setup.py".

The main area is titled "kaggle-api / KaggleSwagger.yaml" and shows the content of the selected file. The file is a YAML document with the following content:

```
1  # Copyright 2018 Kaggle Inc
2  #
3  # Licensed under the Apache License, Version 2.0 (the "License");
4  # you may not use this file except in compliance with the License.
5  # You may obtain a copy of the License at
6  #
7  #     http://www.apache.org/licenses/LICENSE-2.0
8  #
9  # Unless required by applicable law or agreed to in writing, software
10 # distributed under the License is distributed on an "AS IS" BASIS,
11 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12 # See the License for the specific language governing permissions and
13 # limitations under the License.
14 swagger: "2.0"
15 info:
16   description: API for kaggle.com
17   version: "1"
18   title: Kaggle API
19   termsOfService: https://www.kaggle.com/terms
20   contact:
21     name: Contact Kaggle
22     url: https://www.kaggle.com/contact
23   host: www.kaggle.com
24   basePath: /api/v1
25   schemes:
26     - https
27   securityDefinitions:
28     basicAuth:
29       type: basic
```

<https://editor.swagger.io/>

The screenshot shows the Swagger Editor interface. On the left, the API definition for Kaggle is displayed in YAML format. On the right, a list of generated API endpoints for the Kaggle API is shown, each with a method (GET or POST), endpoint path, and a brief description.

```
1 # Copyright 2018 Kaggle Inc
2 #
3 # Licensed under the Apache License, Version 2.0 (the "License");
4 # you may not use this file except in compliance with the License.
5 # You may obtain a copy of the License at
6 #
7 #     http://www.apache.org/licenses/LICENSE-2.0
8 #
9 # Unless required by applicable law or agreed to in writing, software
10 # distributed under the License is distributed on an "AS IS" BASIS,
11 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12 # See the License for the specific language governing permissions and
13 # limitations under the License.
14 swagger: "2.0"
15 info:
16   description: API for kaggle.com
17   version: "1"
18   title: Kaggle API
19   termsOfService: https://www.kaggle.com/terms
20   contact:
21     name: Contact Kaggle
22     url: https://www.kaggle.com/contact
23   host: www.kaggle.com
24   basePath: /api/v1
25   schemes:
26     - https
27   securityDefinitions:
28     basicAuth:
29       type: basic
30     security:
31       - basicAuth: []
32   externalDocs:
33     description: Further Kaggle documentation
34     url: www.kaggle.com
35   tags:
36     - name: kaggle
37   parameters:
38     guidParam:
39       in: path
40       name: guid
41       required: true
42       type: string
43       minimum: 0
44       description: guid specifies location where submission should be
          uploaded
```

kaggle		
GET	/competitions/list	List competitions
GET	/competitions/submissions/list/{id}	List competition submissions
POST	/competitions/{id}/submissions/url/{contentLength}/{lastModifiedDateUtc}	Generate competition submission URL
POST	/competitions/submissions/upload/{guid}/{contentLength}/{lastModifiedDateUtc}	Upload competition submission file
POST	/competitions/submissions/submit/{id}	Submit to competition
GET	/competitions/data/list/{id}	List competition data files
GET	/competitions/data/download/{id}/{fileName}	Download competition data file
GET	/competitions/data/download-all/{id}	Download all competition data files
GET	/competitions/{id}/leaderboard/download	Download competition leaderboard
GET	/competitions/{id}/leaderboard/view	View competition leaderboard
GET	/datasets/list	List datasets

Which eventually permitted the collection of data , such as columns information regarding the dataset for our study.

```
url_kaggle2='https://www.kaggle.com/api/v1/datasets/list/ratman/speech-accent-archive'

response_json = requests.get(url_kaggle2, auth=HTTPBasicAuth(username=os.getenv('KAGGLE_USERNAME'), password=os.getenv('KAGGLE_KEY'))).json()

for file in response_json['datasetFiles']:
    print(file['name'], "number of columns: ", len(file['columns']), "\n")
    for col in file['columns']:
        print(col['name'], ':', col['type'])

reading-passage.txt number of columns:  0
speakers_all.csv number of columns:  11

age : Uuid
age_onset : Uuid
birthplace : String
filename : String
native_language : String
sex : String
speakerid : Uuid
country : String
file_missing? : Boolean
: Unspecified
: Unspecified
```

SPEECH ACCENT ARCHIVE (*Web Scraping*)

Data gathered from this website has been captured through **web scraping**.

Here is an illustration where the list of languages is collected this way

Capturing list of languages:

```
import pandas as pd
from bs4 import BeautifulSoup

languages = []
url1 = "https://accent.gmu.edu/browse_language.php"

response = requests.get(url1)
response.encoding = response.apparent_encoding

if response.status_code == 200:
    soup = BeautifulSoup(response.content, "html.parser")
    ul_languages = soup.find("ul", class_="languagelist")

    for language in ul_languages.find_all("li"):
        languages.append(language.text)

else:
    print("ERROR:", response.status_code)

df_languages = pd.DataFrame(languages, columns=['languages'])
df_languages.to_csv('../data/languages.csv', index=False)
```

DATA CLEANING & EXPLORATORY DATA ANALYSIS

DATA CLEANING

Documentation in progress.

EXPLORATORY DATA ANALYSIS

Documentation in progress...

DATABASE TYPE SELECTION

RELATIONAL DATABASES

Relational Database:

Pros:

- Data Integrity: Relational databases enforce data integrity through constraints like primary keys, foreign keys, and unique constraints, ensuring accurate and consistent data.
- Structured Query Language (SQL): They use SQL for querying and manipulating data, which is a widely understood and standardized language.
- ACID Transactions: Relational databases typically support ACID (Atomicity, Consistency, Isolation, Durability) transactions, ensuring data integrity even in the face of concurrent access.
- Complex Queries: Relational databases are well-suited for complex queries involving multiple tables and relationships.

Cons:

- Scalability: Scaling relational databases can be challenging, especially when dealing with large volumes of data or high transaction rates.
- Schema Changes: Making changes to the database schema can be cumbersome and may require downtime or affect performance.
- Performance: Performance can degrade as the size of the database grows, especially for complex queries involving joins across multiple tables.
- Not Ideal for Unstructured Data: Relational databases are not well-suited for handling unstructured or semi-structured data, such as documents or JSON data.

NON-RELATIONAL DATABASES

Non-Relational Database:

Pros:

- Scalability: Non-relational databases are often designed for horizontal scalability, making them well-suited for handling large volumes of data and high traffic loads.

- Schema Flexibility: Non-relational databases typically offer schema flexibility, allowing for the storage of semi-structured and unstructured data without a fixed schema.
- Performance: Non-relational databases can offer high performance, especially for simple read and write operations, due to their optimized data storage and retrieval mechanisms.
- Handling Big Data: They are well-suited for handling big data and real-time analytics applications, where rapid ingestion and processing of large volumes of data are required.

Cons:

- Data Consistency: Non-relational databases may sacrifice strong consistency guarantees in favor of performance and scalability, leading to eventual consistency models.
- Querying: Querying in non-relational databases may be less expressive compared to SQL in relational databases, especially for complex queries involving joins and aggregations.
- Less Mature Ecosystem: Non-relational databases often have a less mature ecosystem compared to relational databases, with fewer tools, libraries, and community support.
- Learning Curve: Developers accustomed to relational databases may face a learning curve when working with non-relational databases due to differences in data modeling and query languages.

DATABASE CREATION

SQL QUERIES

```
-- Drop SCHEMA if they exist
DROP SCHEMA IF EXISTS accents;

-- Create SCHEMA
CREATE SCHEMA accents;

-- Create database
CREATE DATABASE IF NOT EXISTS accents;

USE accents;

-- Drop tables if they exist
DROP TABLE IF EXISTS kaggle_recs;
DROP TABLE IF EXISTS kaggle_all;
DROP TABLE IF EXISTS saa_lang;
DROP TABLE IF EXISTS saa_rec;
DROP TABLE IF EXISTS saa_locations;
DROP TABLE IF EXISTS saa_bios;

-- Create the table saa_lang
CREATE TABLE saa_lang (
    language_id SMALLINT UNSIGNED NOT NULL AUTO_INCREMENT,
    language VARCHAR(255) NOT NULL,
    PRIMARY KEY (language_id)
);

-- Create the table saa_locations
CREATE TABLE saa_locations (
    location_id SMALLINT UNSIGNED NOT NULL AUTO_INCREMENT,
    location VARCHAR(255) NOT NULL,
    PRIMARY KEY (location_id)
);

-- Create the table saa_rec
CREATE TABLE saa_rec (
    speaker_id SMALLINT UNSIGNED NOT NULL AUTO_INCREMENT,
    filename VARCHAR(255),
    PRIMARY KEY (speaker_id)
);

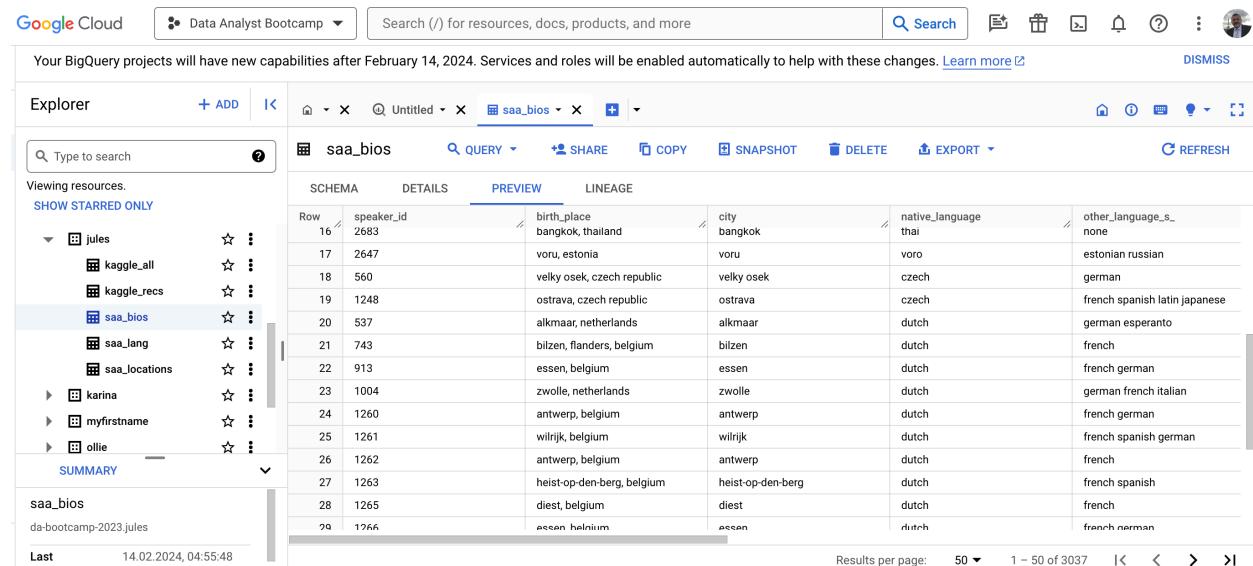
-- Create kaggle_all table
CREATE TABLE kaggle_all (
    age INT NOT NULL,
    age_onset INT NOT NULL,
    birthplace VARCHAR(255),
    filename VARCHAR(255),
    native_language VARCHAR(255),
    sex VARCHAR(255),
    speaker_id SMALLINT UNSIGNED NOT NULL AUTO_INCREMENT,
    country VARCHAR(255),
    file_missing BOOLEAN,
    PRIMARY KEY (speaker_id),
    CONSTRAINT fk_saa_rec_kaggle_all FOREIGN KEY (speaker_id) REFERENCES saa_rec(speaker_id) ON DELETE RESTRICT ON UPDATE CASCADE
);
```

```

-- Create the table kaggle_recs
CREATE TABLE kaggle_recs (
    speaker_id SMALLINT UNSIGNED NOT NULL AUTO_INCREMENT,
    filename VARCHAR(255),
    PRIMARY KEY (speaker_id),
    CONSTRAINT fk_kaggle_all_kaggle_recs FOREIGN KEY (speaker_id) REFERENCES kaggle_all(speaker_id) ON DELETE RESTRICT ON UPDATE CASCADE
);
-- Create the table saa_bios
CREATE TABLE saa_bios (
    speaker_id SMALLINT UNSIGNED NOT NULL AUTO_INCREMENT,
    birth_place VARCHAR(255),
    city VARCHAR(255),
    native_language VARCHAR(255),
    other_languages VARCHAR(255),
    age FLOAT,
    sex VARCHAR(255),
    age_of_english_onset VARCHAR(255),
    english_learning_method VARCHAR(255),
    english_residence VARCHAR(255),
    length_of_english_residence VARCHAR(255),
    language VARCHAR(255),
    language_id SMALLINT UNSIGNED NOT NULL,
    location_id SMALLINT UNSIGNED NOT NULL,
    PRIMARY KEY (speaker_id),
    KEY idx_fk_location_id (location_id),
    KEY idx_fk_language_id (language_id),
    CONSTRAINT fk_saa_bios_location_id FOREIGN KEY (location_id) REFERENCES saa_locations (location_id) ON DELETE RESTRICT ON UPDATE CASCADE,
    CONSTRAINT fk_saa_bios_language_id FOREIGN KEY (language_id) REFERENCES saa_lang (language_id) ON DELETE RESTRICT ON UPDATE CASCADE
);

```

Cloud: BigQuery



The screenshot shows the Google Cloud BigQuery interface. At the top, there is a navigation bar with 'Google Cloud' and a dropdown for 'Data Analyst Bootcamp'. A search bar contains the placeholder 'Search (/) for resources, docs, products, and more'. On the far right of the top bar are several icons: a profile picture, a bell, a question mark, and a 'DISMISS' button.

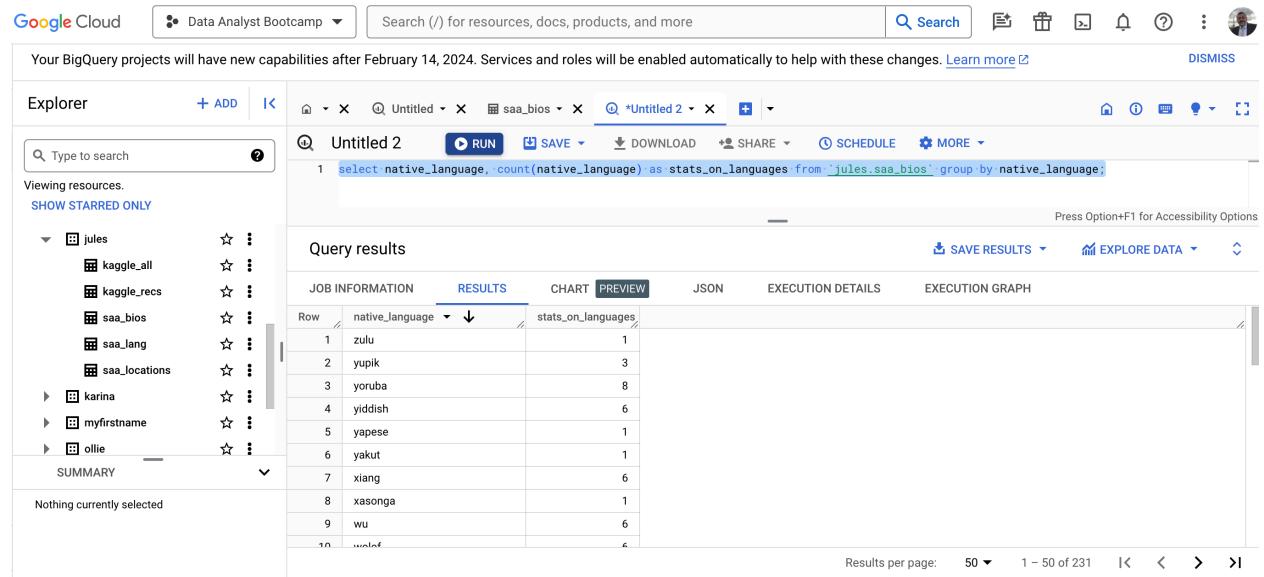
Below the top bar, a message states: 'Your BigQuery projects will have new capabilities after February 14, 2024. Services and roles will be enabled automatically to help with these changes. [Learn more](#)' with a 'DISMISS' button.

The main area is titled 'Explorer' and shows a tree view of datasets and tables. Under the 'jules' dataset, the 'saa_bios' table is selected. The table has columns: speaker_id, birth_place, city, native_language, and other_languages. The preview pane shows 29 rows of data, including entries for Bangkok, Thailand; Voru, Estonia; and various cities in the Czech Republic, Germany, Belgium, and the Netherlands.

Row	speaker_id	birth_place	city	native_language	other_languages
16	2683	bangkok, thailand	bangkok	thai	none
17	2647	voru, estonia	voru	voru	estonian russian
18	560	velky osek, czech republic	velky osek	czech	german
19	1248	ostrava, czech republic	ostrava	czech	french spanish latin japanese
20	537	alkmaar, netherlands	alkmaar	dutch	german esperanto
21	743	bilzen, flanders, belgium	bilzen	dutch	french
22	913	essen, belgium	essen	dutch	french german
23	1004	zwolle, netherlands	zwolle	dutch	german french italian
24	1260	antwerp, belgium	antwerp	dutch	french german
25	1261	wilrijk, belgium	wilrijk	dutch	french spanish german
26	1262	antwerp, belgium	antwerp	dutch	french
27	1263	heist-op-den-berg, belgium	heist-op-den-berg	dutch	french spanish
28	1265	diest, belgium	diest	dutch	french
29	1266	essen, belgium	essen	dutch	french german

At the bottom left, it says 'Last 14.02.2024, 04:55:48'. On the bottom right, there are buttons for 'Results per page: 50', '1 - 50 of 3037', and navigation arrows.

```
select native_language, count(native_language) as stats_on_languages from `jules.saa_bios` group by native_language;
```



The screenshot shows the Google Cloud BigQuery interface. At the top, there's a navigation bar with 'Google Cloud' and a dropdown for 'Data Analyst Bootcamp'. A search bar contains 'Search (/) for resources, docs, products, and more' with a 'Search' button. Below the search bar, a message says 'Your BigQuery projects will have new capabilities after February 14, 2024. Services and roles will be enabled automatically to help with these changes.' with a 'Learn more' link and a 'DISMISS' button.

The main area has tabs for 'Explorer' and 'Untitled 2'. Under 'Untitled 2', there's a 'RUN' button, a 'SAVE' dropdown, a 'DOWNLOAD' button, a 'SHARE' dropdown, a 'SCHEDULE' button, and a 'MORE' dropdown. The message 'Press Option+F1 for Accessibility Options' is at the top right.

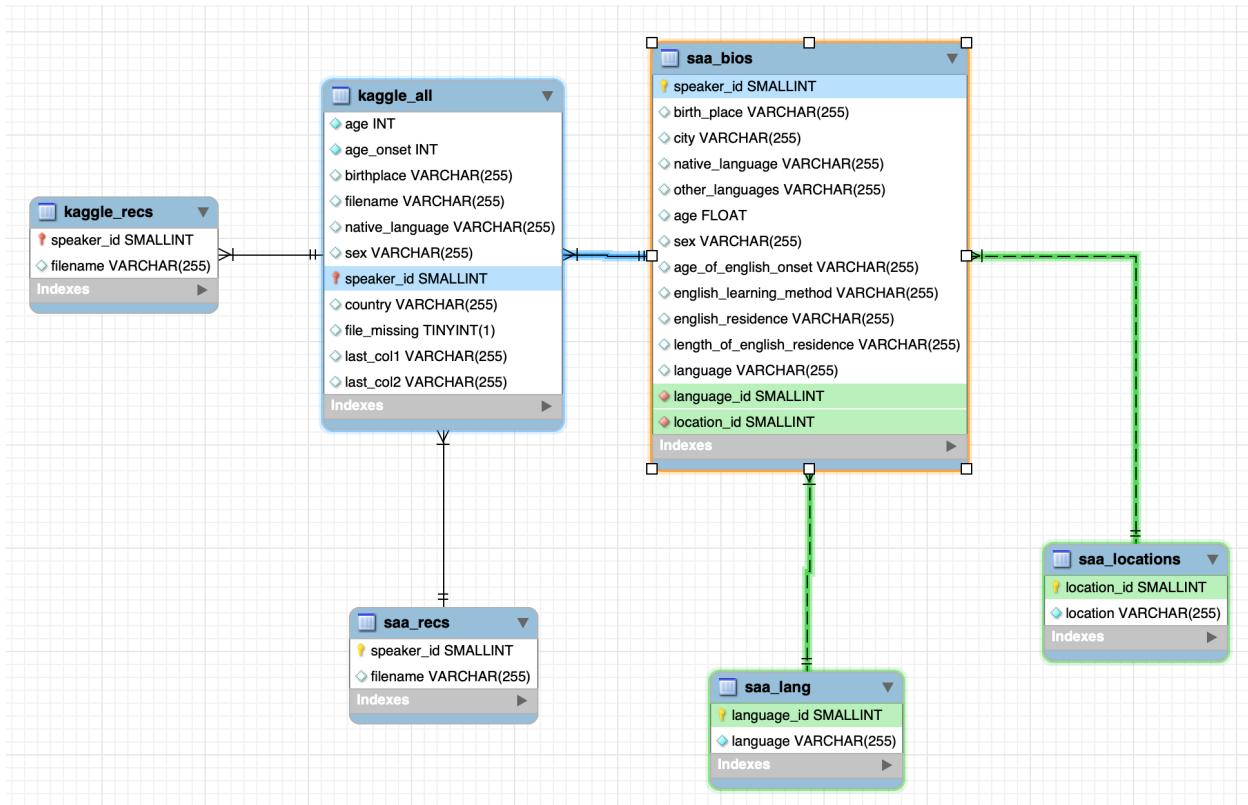
The 'Query results' section has tabs for 'RESULTS' (which is selected), 'CHART', 'PREVIEW', 'JSON', 'EXECUTION DETAILS', and 'EXECUTION GRAPH'. The results table has columns 'Row', 'native_language', and 'stats_on_languages'. The data is as follows:

Row	native_language	stats_on_languages
1	zulu	1
2	yupik	3
3	yoruba	8
4	yiddish	6
5	yapese	1
6	yakut	1
7	xiang	6
8	xasonga	1
9	wu	6
10	walef	6

At the bottom, it says 'Results per page: 50 ▾ 1 – 50 of 231 |< < > >|'.

ENTITY RELATIONSHIP DIAGRAM

Here is the ERD of our Schema.



GENERAL DATA PROTECTION REGULATION (GDPR)

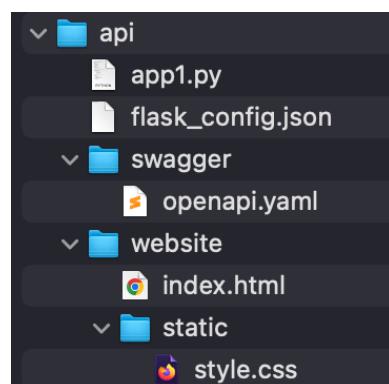
TYPE OF DATA

The usage of soundtracks used in this project is compliant with GDPR norms.

API DEVELOPMENT

CREATING REST API WITH FLASK

FOLDERS STRUCTURE



YAML FILE

openapi.yaml.

```
openapi: 3.0.0
info:
  title: Swagger
  version: 1.0.0
paths:
  /speaker_info/{speaker_id}:
    get:
      summary: Get speaker information
      parameters:
        - in: path
          name: speaker_id
          required: true
          schema:
            type: integer
      responses:
        '200':
          description: Successful operation
          content:
            application/json:
              schema:
                type: object
                properties:
                  native_language:
                    type: string
                  age:
                    type: number
                  sex:
                    type: string
        '404':
          description: Speaker not found
  /locations:
    get:
      summary: Get all locations
      responses:
        '200':
          description: Successful operation
          content:
            application/json:
              schema:
                type: array
                items:
                  type: string
```

APP1.PY

```
from flask import Flask, jsonify, request
from flask_swagger_ui import get_swaggerui_blueprint
import mysql.connector
import pymysql
import os

app = Flask(__name__)

db = mysql.connector.connect(
    host="localhost",
    user="root",
    password=os.getenv('MySQLpwd'),
    database="accents"
)
cursor = db.cursor()

SWAGGER_URL = '/swagger'
API_URL = '/static/openapi.yaml'
swaggerui_blueprint = get_swaggerui_blueprint(
    SWAGGER_URL,
    API_URL,
    config={'app_name': "Swagger"}
)
app.register_blueprint(swaggerui_blueprint, url_prefix=SWAGGER_URL)

@app.route('/speaker_info/<int:speaker_id>')
def get_speaker_info(speaker_id):
    try:
        query = "SELECT native_language, age, sex FROM saa_bios WHERE speaker_id = %s"
        cursor.execute(query, (speaker_id,))
        result = cursor.fetchone()
        if result:
            speaker_info = {
                'native_language': result[0],
                'age': result[1],
                'sex': result[2]
            }
            return jsonify(speaker_info)
        else:
            return jsonify({'message': 'Speaker not found'}), 404
    except Exception as e:
        return jsonify({'message': str(e)}), 500

@app.route('/locations')
def get_locations():
    try:
        query = "SELECT location FROM saa_locations"
        cursor.execute(query)
        results = cursor.fetchall()
        locations = [result[0] for result in results]
        return jsonify(locations)
    except Exception as e:
        return jsonify({'message': str(e)}), 500

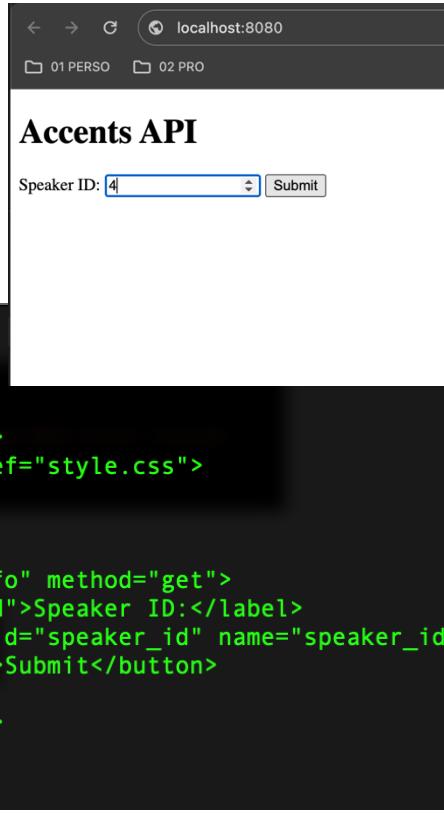
if __name__ == '__main__':
    app.run(port=8080)
```

```
~/projects/accent_recognition/api$ flask --app app1 run --port 8080
 * Serving Flask app 'app1'
 * Debug mode: off
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
 * Running on http://127.0.0.1:8080
Press CTRL+C to quit
```

WEBSITE

HTML

```
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <title>Accents API</title>
    <link rel="stylesheet" href="style.css">
</head>
<body>
    <h1>Accents API</h1>
    <form action="/speaker_info" method="get">
        <label for="speaker_id">Speaker ID:</label>
        <input type="number" id="speaker_id" name="speaker_id">
        <button type="submit">Submit</button>
    </form>
    <div id="locations"></div>
</body>
</html>
```



CSS

```
body {
    font-family: Arial, sans-serif;
    margin: 20px;
}

h1 {
    color: #333;
}

form {
    margin-bottom: 20px;
}

label {
    font-weight: bold;
}

input[type="number"] {
    width: 100px;
    margin-right: 10px;
}

button {
    padding: 5px 10px;
    background-color: #007bff;
    color: #fff;
    border: none;
    cursor: pointer;
}

button:hover {
    background-color: #0056b3;
}

#locations {
    margin-top: 20px;
}
```

MACHINE LEARNING

MODELS SELECTION

XXX.

MODELS APPLICATION

Features and Target used in the following models

```
: features, target = num_var, 'native_language'  
if target in features: features.remove(target)
```

Scaling & Splitting dataset

```
X_train_scaled, X_test_scaled, y_train, y_test = jt.scaling(df, features, target,StandardScaler())
```

Models application: in progress...(to be defined, but likely **CNN** or **RNN**)

Evaluation Report:

```
: for model in models:  
    y_pred = jt.model_fit_predict(model,X_train_scaled, X_test_scaled, y_train, y_test)  
    jt.eval_report(y_test,y_pred)
```