

# Master in Data Science AUEB 2024-2025

## Practical Data Science - Assignment 2

Stylianos Giagkos f3352410

## Food Hazard Detection Challenge

This repository contains the solution for the **Food Hazard Detection Challenge**. The challenge involves classifying food safety-related incidents based on short titles and long descriptions. The solution leverages both **Finetuned PubMedBERT** and **LightGBM (LGBM)** models for classification into hazard-category, product-category, hazard, and product classifications.

### Overview

#### Project Workflow:

1. **Exploratory Data Analysis (EDA)**: Data cleaning, exploration, and visualization.
2. **Modeling**: Two different approaches are used for classification:
  - **Finetuned PubMedBERT**: Fine-tuned on the dataset to classify food hazard-related texts.
  - **LightGBM**: A gradient boosting model for classification based on features derived from the data.
3. **Evaluation**: Both models are evaluated on performance metrics such as accuracy, precision, recall, and F1-score.
4. **Training and Submission**: Generation of final predictions based on the final optimal model and submission in the required format on CodaLab competition.

#### Subtasks (Performed Separately for Title and Text):

##### Subtask 1:

- **Classify hazard-category**: Classifies the general hazard type.
- **Classify product-category**: Classifies the general product type.

##### Subtask 2:

- **Classify hazard**: Classifies the specific hazard type.
- **Classify product**: Classifies the specific product type.

### Repository Files Description:

The following files are included in the repository:

## Benchmarks Notebooks

- **[Augmented Train Set Benchmark Models Finetuned PubMedBERT PDS A2.ipynb]**  
Notebook for evaluating benchmark models using an **augmented training set** with PubMedBERT.
- **[Benchmark Models Finetuned PubMedBERT PDS A2 Food Hazard Detection.ipynb]**  
Notebook for running and evaluating benchmark models using PubMedBERT with the **initial training set**.
- **[Benchmark Models LGBM PDS A2 Food Hazard Detection.ipynb]**  
Notebook for experimenting with LightGBM models for food hazard detection using the **initial training set**.

## Submission Notebooks

- **[Augmented Submission Model Finetuned PubMedBERT PDS A2.ipynb]**  
Submission notebook for a model finetuned on an **augmented training set** using PubMedBERT.
- **[Submission Model Finetuned PubMedBERT PDS A2 Food Hazard Detection.ipynb]**  
Submission notebook for a model finetuned on the **initial training set** using PubMedBERT.
- **[Submission Model LGBM PDS A2 Food Hazard Detection.ipynb]**  
Submission notebook for LightGBM models trained on the **initial training set**.

## Extra Notebooks

- **[EDA Notebook PDS A2 Food Hazard Detection.ipynb]**  
Notebook for performing Exploratory Data Analysis (EDA) on the **initial training set**.
- **[Incidents Augmentation using Techniques.ipynb]**  
Notebook for applying data augmentation techniques to incident descriptions and titles, creating an **augmented training set** for imbalanced classes especially in categories of **hazard** and **product**.

## Requirements

Make sure to install the required dependencies before running the code. You can use the following pip command to install the necessary packages:

```
pip install torch lightgbm pandas scikit-learn matplotlib tqdm transformers nltk numpy
```

### Additional Libraries:

- pandas
- re
- nltk
- scikit-learn
- torch
- transformers
- lightgbm
- matplotlib
- numpy

## How to Re-run the Solution

### Step 1: Data Preparation

Ensure that the dataset is correctly placed in the expected directory. Adjust file paths if necessary.

### Step 2: Exploratory Data Analysis

Run EDA PDS A2 Food Hazard Detection.ipynb to clean and visualize the data.

### Step 3: Model Training

You can choose between:

- **Finetuned PubMedBERT:** Run BENCHMARKS Finetuned PubMedBERT PDS A2 Food Hazard Detection.ipynb.
- **LightGBM:** Run BENCHMARKS LGBM PDS A2 Food Hazard Detection .ipynb.

### Step 4: Model Evaluation

Both models evaluate accuracy, precision, recall, and F1-score.

### Step 5: Retraining for Submission

Retrain using the submission notebooks to create final submissions for competitions based on the best model of the respective benchmarks notebooks:

- **For PubMedBERT:** Run SUBMISSION Finetuned PubMedBERT PDS A2 Food Hazard Detection.ipynb.
- **For LightGBM:** Run SUBMISSION LGBM PDS A2 Food Hazard Detection .ipynb.

## Results

### Finetuned PubMedBERT Model (Title-based):

Task	F1-Score
hazard-category (Title)	0.8288
product-category (Title)	0.7494
hazard (Title)	0.5899
product (Title)	0.2172

### Finetuned PubMedBERT Model (Text-based):

Task	F1-Score
hazard-category (Text)	0.9459
product-category (Text)	0.7583
hazard (Text)	0.8166
product (Text)	0.2331

### LightGBM Model (Title-based):

Task	F1-Score
hazard-category (Title)	0.7614
product-category (Title)	0.5926
hazard (Title)	0.5533
product (Title)	0.0798

### LightGBM Model (Text-based):

Task	F1-Score
hazard-category (Text)	0.9129
product-category (Text)	0.6682
hazard (Text)	0.7671
product (Text)	0.0479

## Competition Results for Subtask 1 (Hazard-category, Product-category)

The following shows the results for hazard-category and product-category tasks for Subtask 1:

Submission File	Score	Status
submission.zip	—	Failed
submission.zip	0.6428057851	Finished
lgbm_submission_title.zip	0.6428057851	Finished
submission_v3.zip	0.6252022003	Finished
submission_finetuned_PubMedBERT.zip	0.6992400644	Finished

## Competition Results for Subtask 2 (Hazard, Product)

The following shows the results for hazard and product tasks for Subtask 1:

*But the rating shows for Sub Task 2 a score of 0.3313 .*

Submission File	Score	Status
submission.zip	—	Failed
submission.zip	0.00000	Finished
lgbm_submission_title.zip	0.00000	Finished
submission_v3.zip	0.00000	Finished
submission_finetuned_PubMedBERT.zip	0.00000	Finished

## Overfitting Indications from Competition Results

The competition results suggest potential overfitting in the models. The scores for the `submission_finetuned_PubMedBERT.zip` file were significantly higher for the training set compared to the test set, which may indicate that the model is overfitting to the training data. To mitigate this, I attempted data augmentation in the **Incidents Augmentation using techniques.ipynb** file. The benchmark results using the augmented training data are as follows:

### Collected F1-Scores for Title-Focused Classification (Augmented Data)

Task	F1-Score
hazard-category	0.9266
product-category	0.8972
hazard	0.8169
product	0.5898

### Collected F1-Scores for Text-Focused Classification (Augmented Data)

Task	F1-Score
hazard-category	0.9643
product-category	0.8834
hazard	0.8887
product	0.6004

The submission file associated with this augmented training approach was **Augmented\_Submission\_Model\_Finetuned\_PubMedBERT\_PDS\_A2\_Food\_Hazard\_Detection**. The leaderboard scores for this submission were:

- **ST1:** 0.6870518521 (submission\_augmented\_train\_set\_finetunedPUBMEDBERT.zip) on 11/20/2024 at 17:27:18
- **ST2:** 0.0 (submission\_finetuned\_PubMedBERT.zip) on 11/20/2024 at 17:41:21

Despite the data augmentation efforts, there was no significant improvement in the competition leaderboard scores. For instance, ST1 achieved a score of **0.6870518521**, and the ST2 submission showed discrepancies, with the leaderboard reflecting an approximate score of **0.35**. This suggests that data augmentation did not substantially address the overfitting issue or enhance generalization performance in the competition context. Further investigation and alternative approaches may be required to improve model performance.

## Explanation for Model Performance Despite Augmented Dataset

I am trying to explain why the augmented dataset did not improve the model's performance in the following points:

1. **Nature of Augmentation:** The augmentation techniques I applied might not have introduced enough diversity or complexity in the data. This could mean the model didn't get new, meaningful examples to improve its ability to generalize.
2. **Overfitting:** Despite using an augmented dataset, the model might still be overfitting to both the original and augmented data. If the augmented data is too similar to the original, the model may memorize patterns rather than learning to generalize, leading to poor performance on unseen data.

## Conclusion

This solution combines advanced NLP (Finetuned PubMedBERT) and traditional machine learning (LightGBM) techniques to classify food hazard data. The solution addresses both general and specific hazard and product classification tasks. Both models were trained and evaluated, with their results summarized above. The submission models are retrained based on the benchmark results to generate the final predictions. You can re-run the solution by following the provided instructions and reproduce the results with the corresponding datasets and models.

## References:

- ChatGPT: For quick assistance, code help, debugging, and guidance on various data science topics.
- Stack Overflow: For troubleshooting and solutions to coding challenges.
- Towards Data Science and Medium: For articles and tutorials on data science techniques and best practices.