



**IMT Atlantique**

Bretagne-Pays de la Loire

École Mines-Télécom

# Analyse des données de santé : épidémiologie et aide à la décision

## PROJET : Effet de la musique sur la santé mentale

BEN JEMAA Yosr, DRIRA Yosr

# SOMMAIRE

## 1. EXPLORATION DES DONNÉES

- 1.1 Contexte
- 1.2 Problématique
- 1.3 Démographie & comportements
- 1.4 Santé mentale
- 1.5 Exploration des relations

## 2. ANALYSES SUPERVISÉES

- 2.1 Régression
- 2.2 Classification

## 3. CONCLUSION



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# I - EXPLORATION DES DONNÉES



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

## 1.1 Contexte

- La thérapie musicale utilise la musique pour améliorer le stress, l'humeur et la santé mentale.
- Le jeu de données cherche à trouver des liens entre les préférences musicales et la santé mentale.
- Les résultats pourraient guider l'application de la thérapie musicale ou offrir des perspectives sur le fonctionnement de l'esprit.



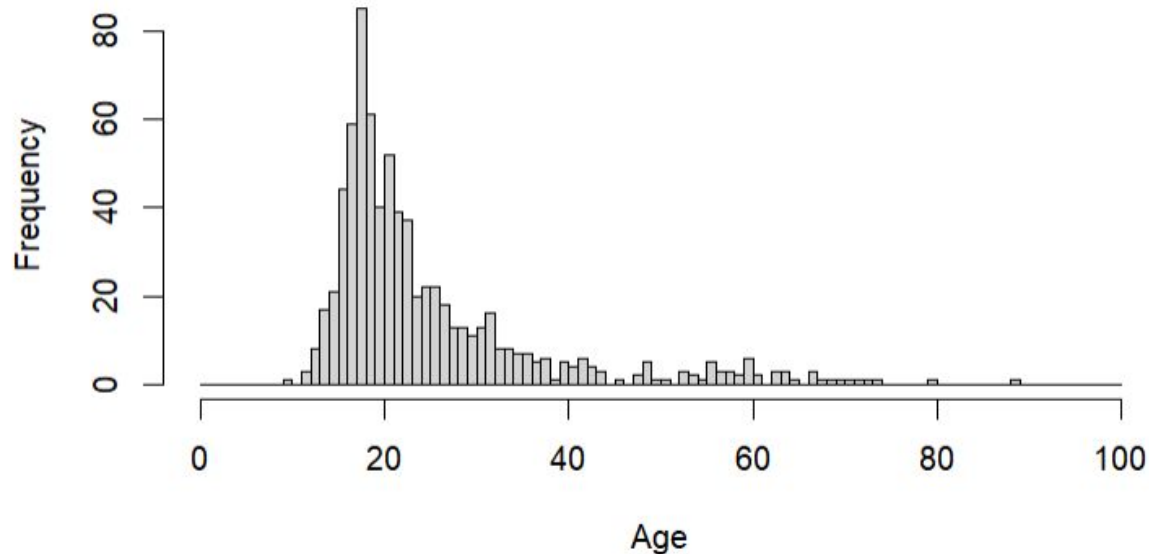
**Jeu de données:** 700 responses (Music x Mental Health Traits) - 33 columns

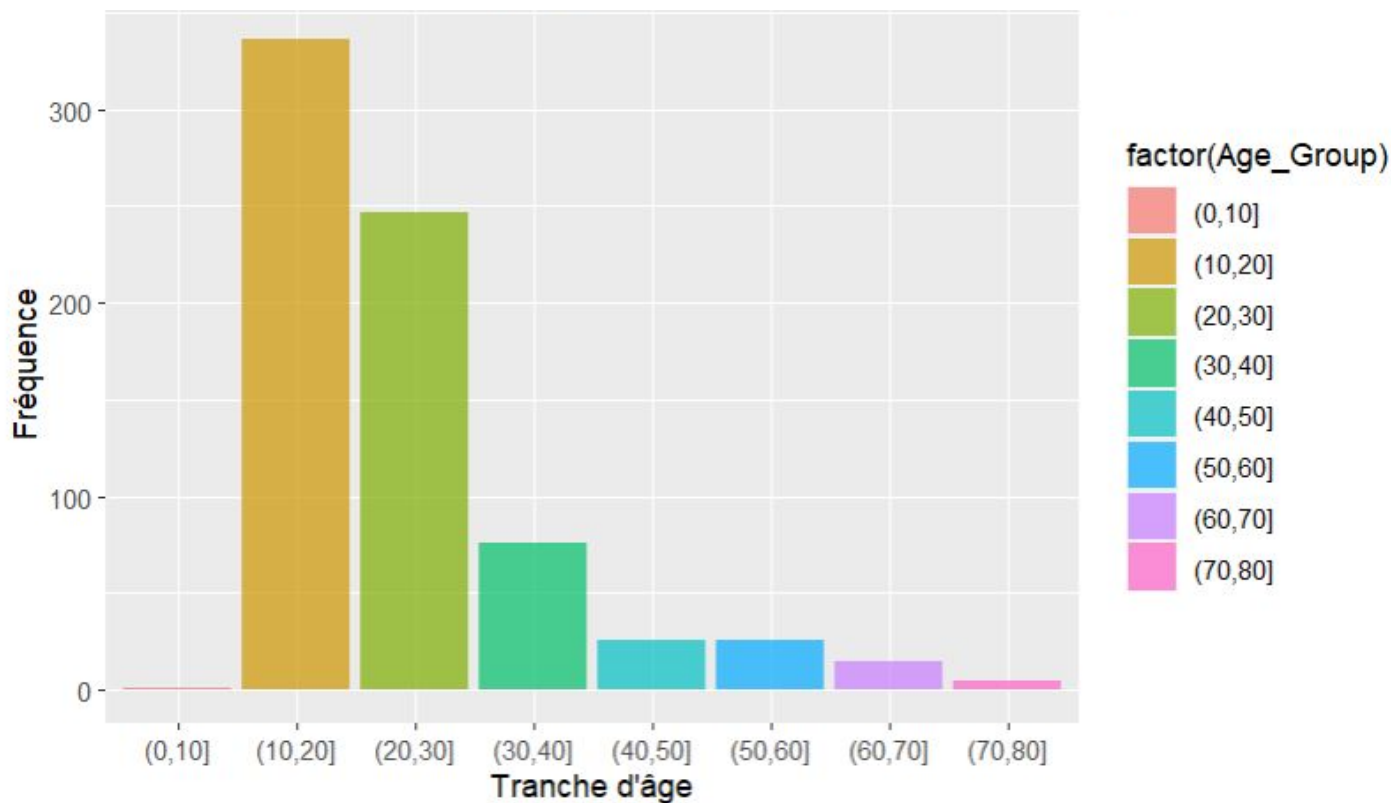
Variables Quantitatives	Variables Qualitatives
Timestamp	Primary streaming service
Age	While working
Hours per day	Instrumentalist
BPM	Composer
Frequency [genre]	Fav genre
Anxiety/ Depression/ Insomnia/ OCD	Exploratory
	Foreign Language
	Music effects



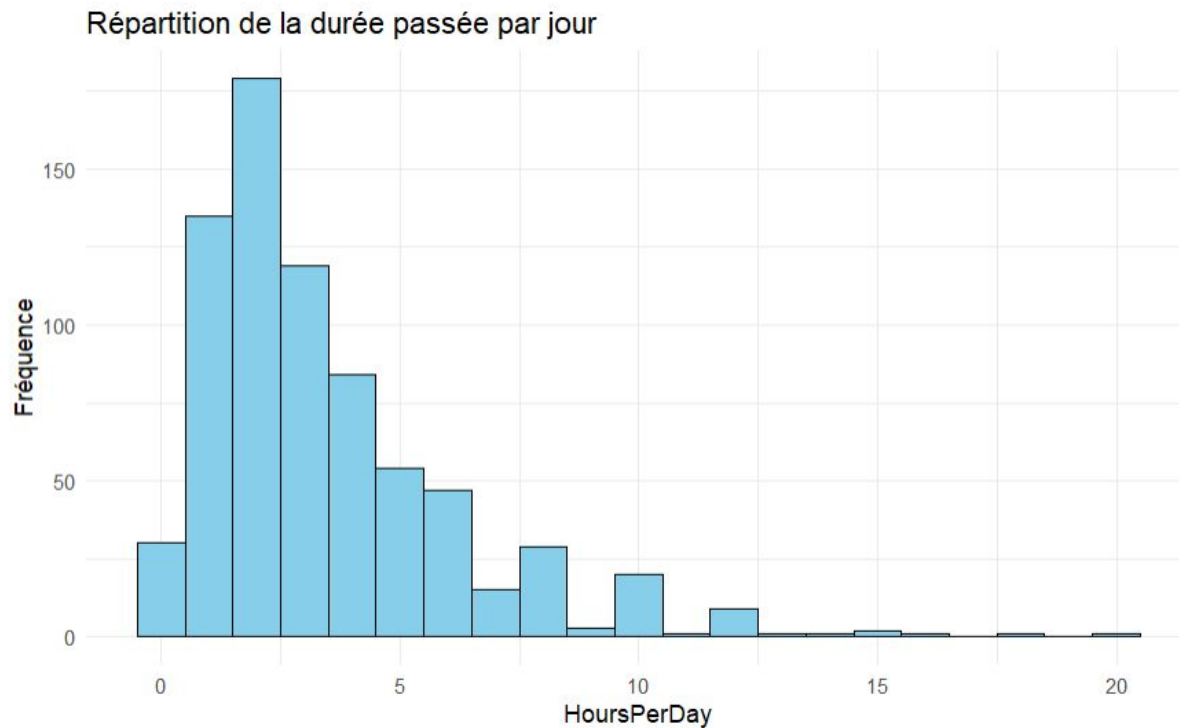
Comment la préférence musicale et les habitudes d'écoute peuvent-elles influencer les niveaux d'anxiété, de dépression, d'insomnie et les troubles obsessionnels-compulsifs (TOC) chez les individus ?

Répartition de l'âge

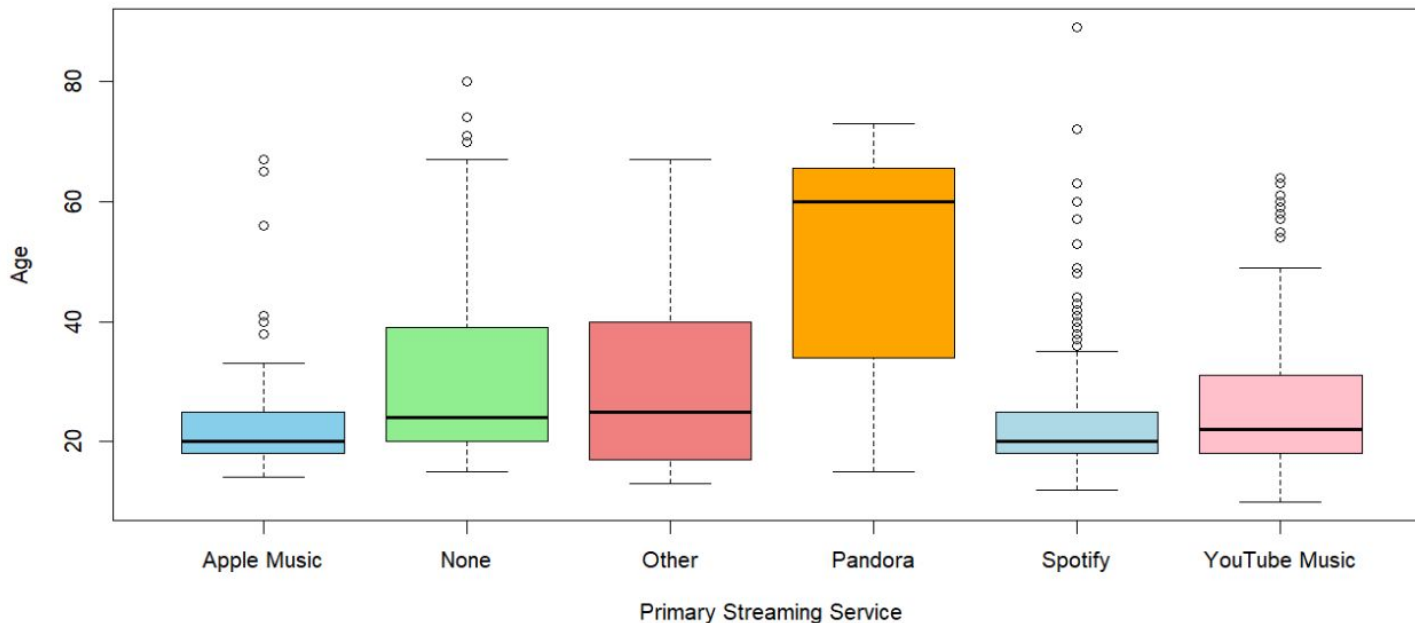




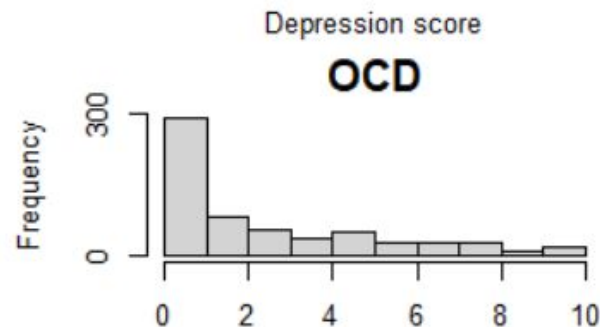
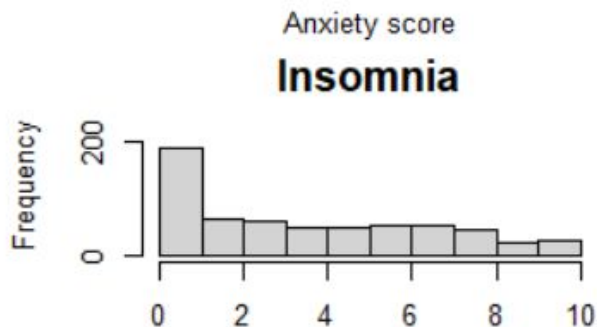
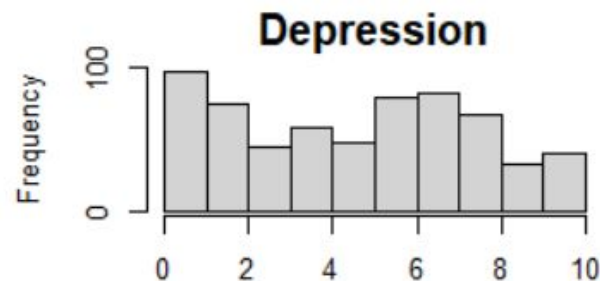
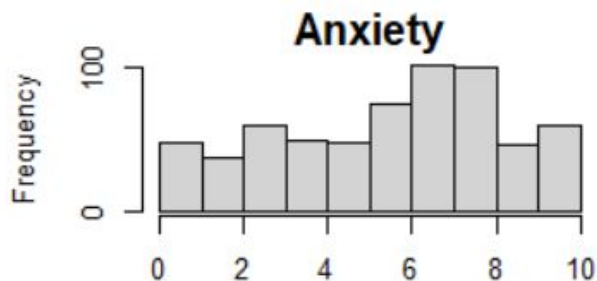




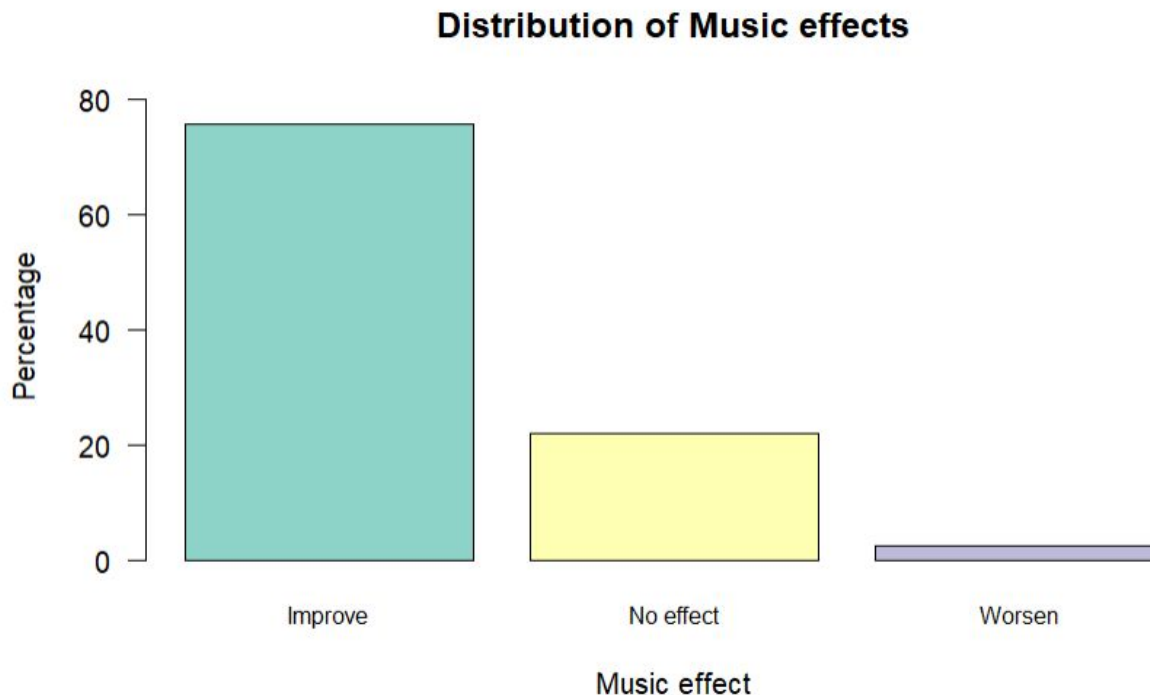
### Diagrammes en BoxPlot de l'âge en fonction du service de streaming



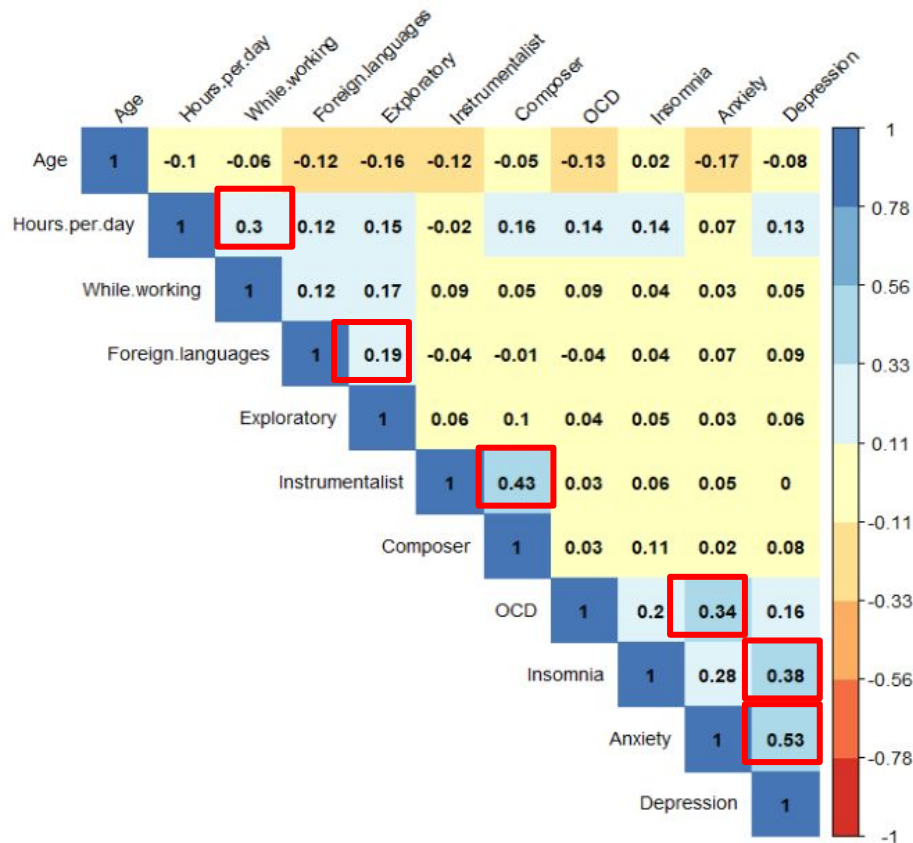
### Histogrammes de répartition de fréquence des degré des indicateurs



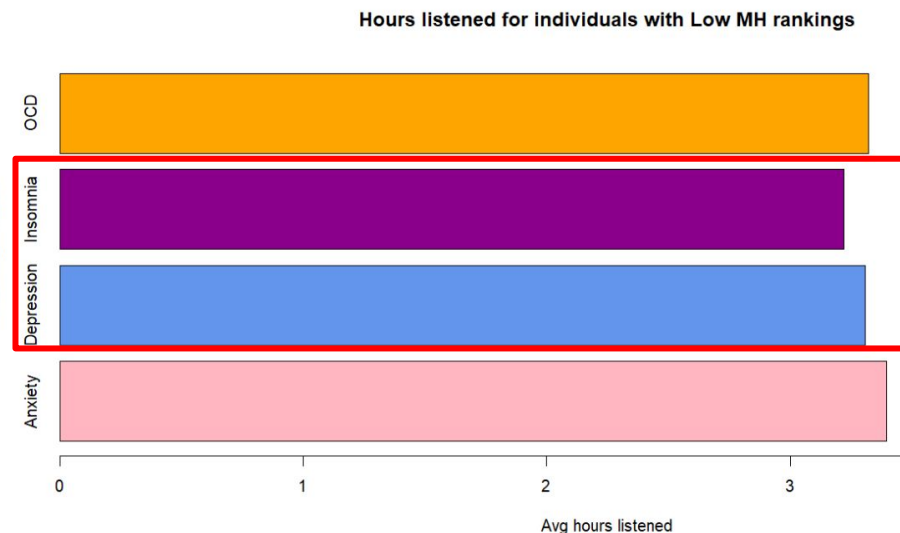
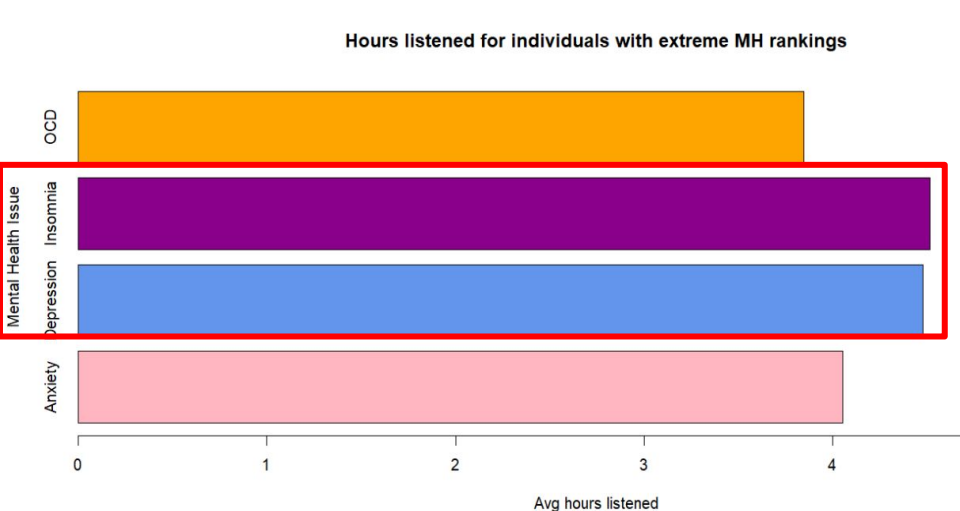
### Distribution de la répartition de l'effet de la musique

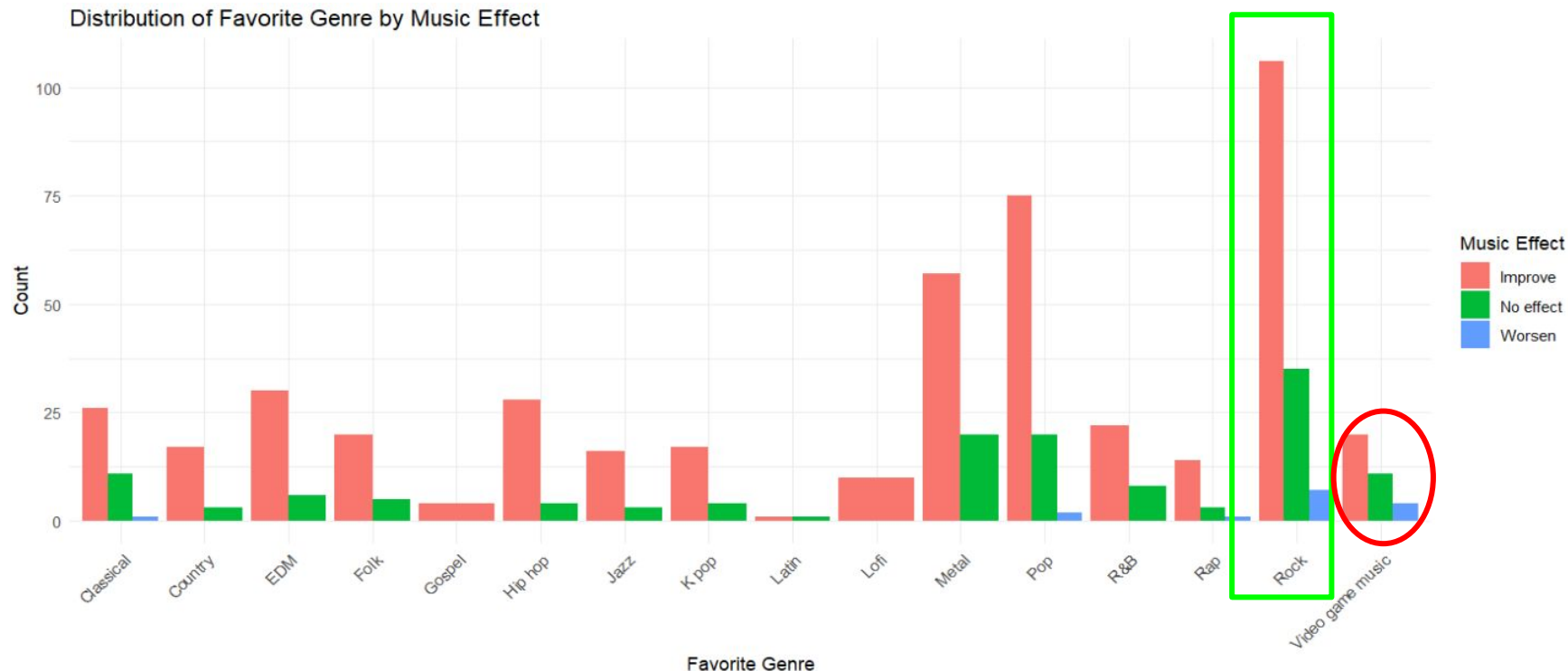


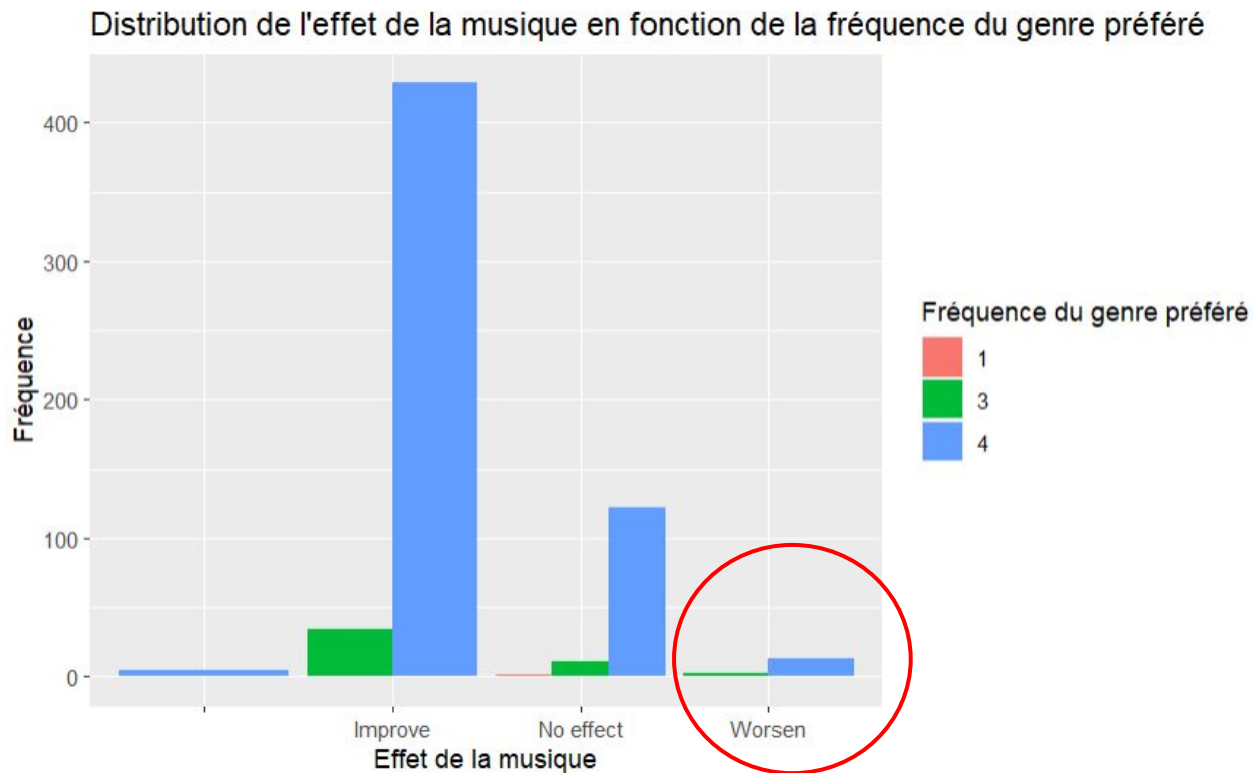
### Heatmap



### Impact du nombre d'heures d'écoute de musique sur les indicateurs

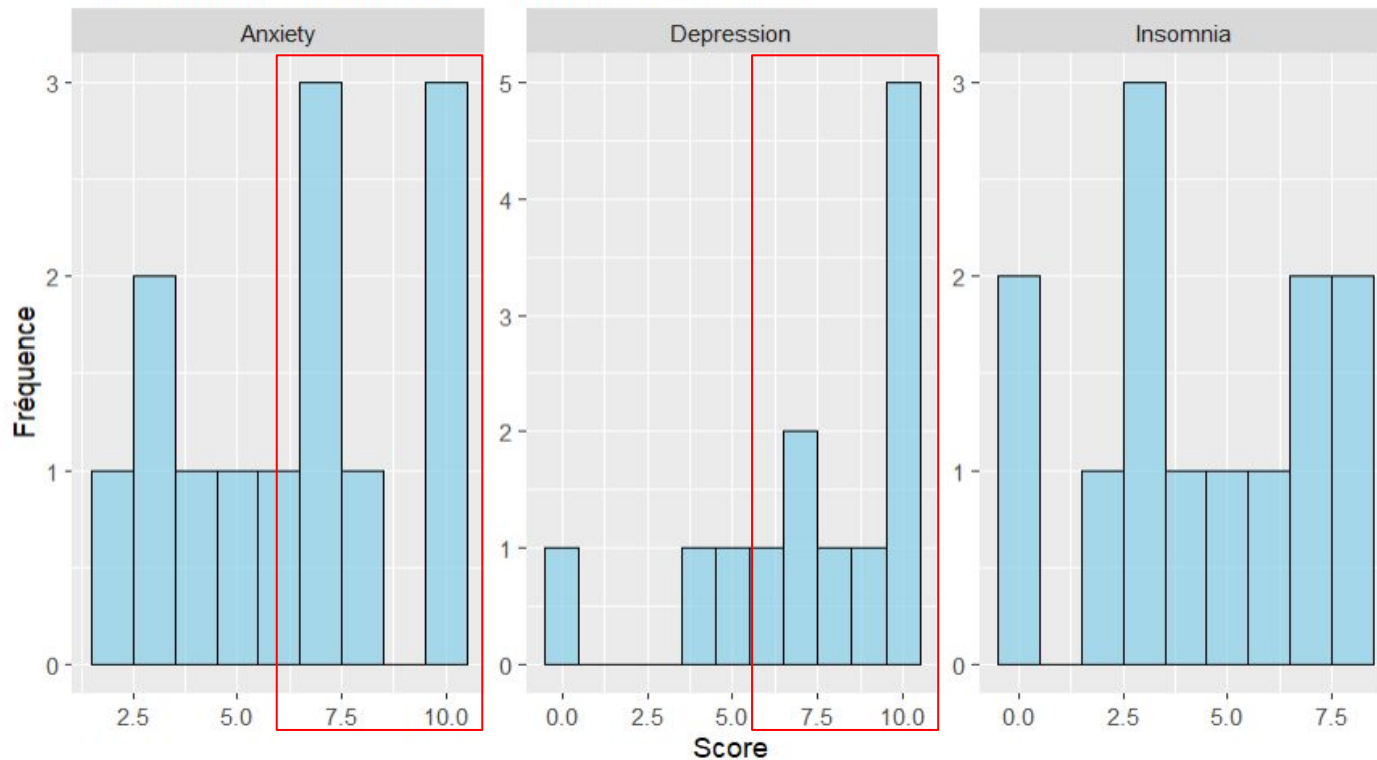


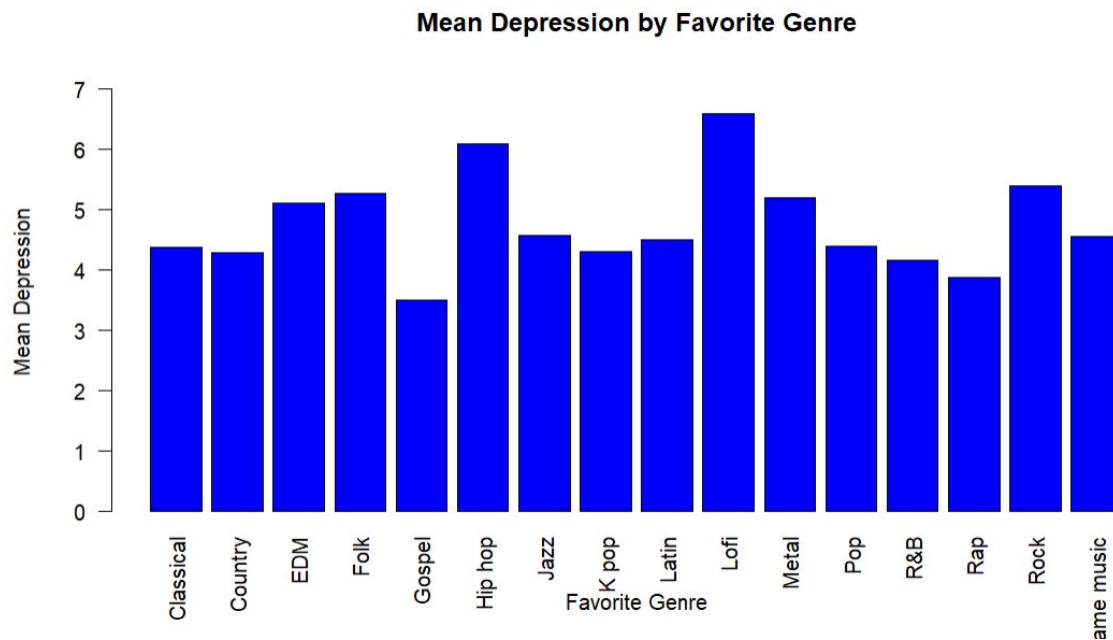




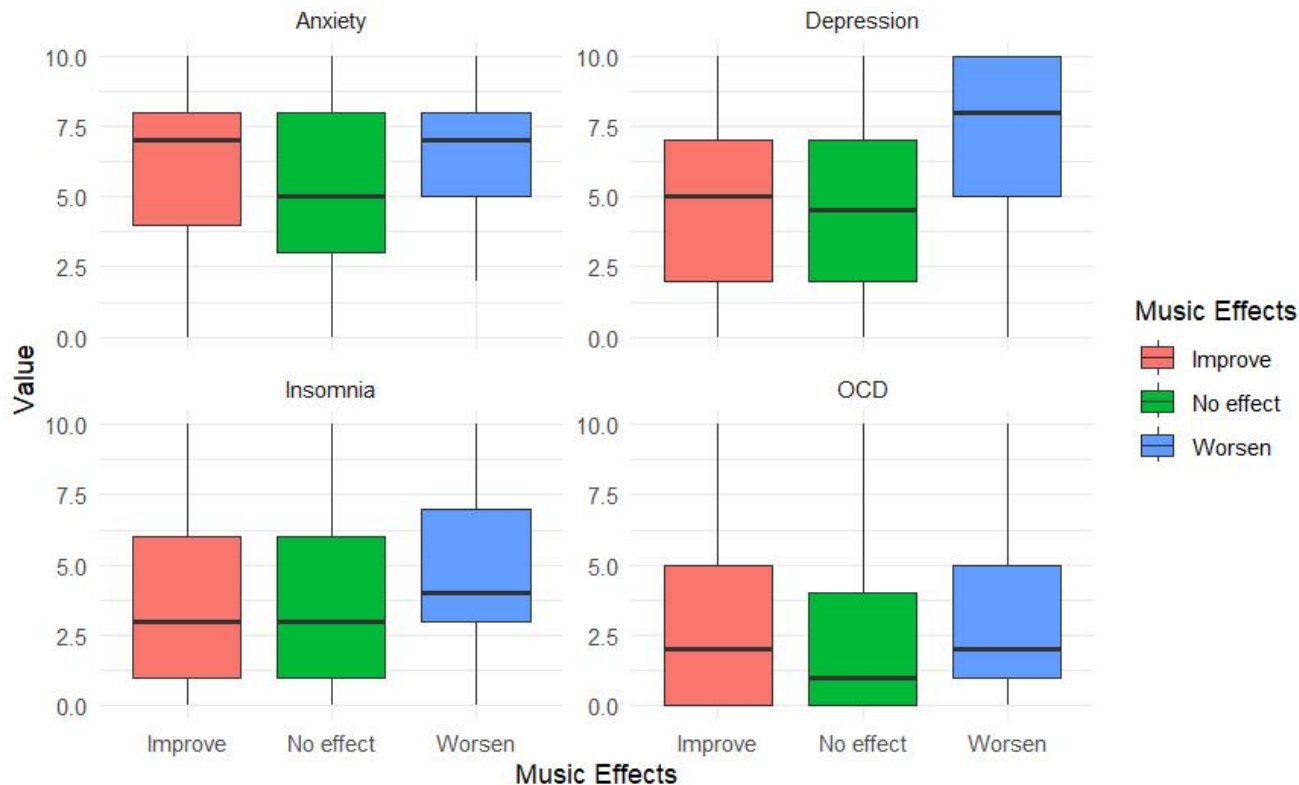


Histogramme de Anxiety, Depression et Insomnia





### Boxplots des mesures de santé mentale selon les effets de la musique

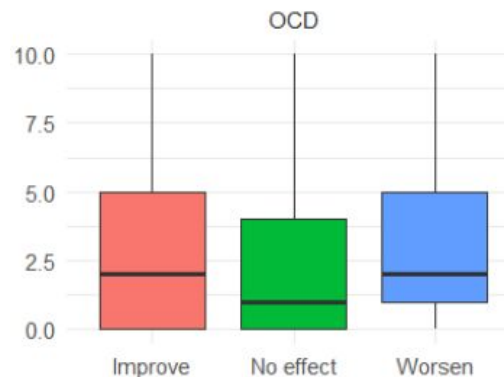


### Relation entre les Indicateurs de santé mentale et effet de la musique

Le test de **Kruskal-Wallis** est une méthode non paramétrique utilisée pour comparer les valeurs médianes de trois groupes indépendants ou plus.

**Hypothèse nulle ( $H_0$ ) :**

L'OCD n'a pas d'impact significatif sur la perception de l'effet de la musique



### Résultats du test Kruskal-Wallis Test

```
data: Music.effects by OCD  
Kruskal-wallis chi-squared = 15.329, df = 12, p-value = 0.2239
```

Kruskal-wallis rank sum test

```
data: Music.effects by Anxiety  
Kruskal-wallis chi-squared = 16.879, df = 11, p-value = 0.1115
```

Kruskal-wallis rank sum test

```
data: Music.effects by Depression  
Kruskal-wallis chi-squared = 30.829, df = 11, p-value = 0.001173
```

→ Le COD / l'anxiété n'ont pas un impact significatif sur la perception de l'effet de la musique.

→ La dépression et l'effet de la musique sont liés.

## II. ANALYSES SUPERVISÉES



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# Régression

**Objectif** : Prédire les scores de santé mentale (MH Rankings) en fonction des caractéristiques disponibles.

## Approches de modélisation :

Régression linéaire

Random Forest

XGBoost.



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

- **Entraînement** : Nous avons entraîné le modèle en utilisant 80 % de l'ensemble de données, en réservant les 20 % restants pour les tests.
- **Évaluation** :

$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

$\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$  are predicted values

$y_1, y_2, \dots, y_n$  are observed values

$n$  is the number of observations

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

Where:

$\hat{y}_i$  = Predicted value for the  $i^{\text{th}}$  data point

$y_i$  = Actual value for the  $i^{\text{th}}$  data point

$n$  = number of observations



## Evaluation des modèles

Target = Dépression

Métrique	Linear Regression	Random Forest	XGBoost
RMSE	3.09	3.08	3.16
MAE	2.62	2.64	2.65

### Méthodes d'Optimisation :

- **Recherche de Grille** : Exploration d'un ensemble prédéfini de hyperparamètres pour trouver la combinaison optimale.
- **Validation Croisée** : Évaluation des performances du modèle sur plusieurs sous-ensembles des données pour estimer sa capacité de généralisation

```
RMSE was used to select the optimal model using the smallest value.  
The final values used for the model were nrounds = 50, max_depth = 3, eta =  
0.1, gamma = 0, colsample_bytree =  
1, min_child_weight = 1 and subsample = 1.  
XGBoost Model Evaluation for Depression:  
RMSE: 3.129314  
MAE: 2.680435
```

⇒ Il n'y a pas une grande différence

## Normalization de la target: Score de dépression

Métrique	Random Forest	XGBoost
RMSE	0.30	0.34
MAE	0.26	0.28

# Classification

## **Objectif :**

Prédire l'effet de la musique à partir des indicateurs (Dépression, Anxiété, Insomnie, TOC/OCD).

## **Approches de modélisation :**

Arbre de décision

Random Forest

Multinomial Logistic Regression



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

## 1. Arbre de décision

**Modèle :** Arbre de décision

**Bibliothèque utilisée:** caret

```
tree_model <- rpart(Music.effects ~ OCD + Anxiety + Depression + Insomnia, data = train_data, method = "class")
```

**Matrice de confusion:**

predicted_music_effect	Improve	No effect	Worsen
Improve	105	31	3
No effect	2	2	0
Worsen	0	0	0

## 1. Arbre de décision

### Métriques d'évaluations:

	'Improve' Class	'No effect' Class	'Worsen' Class
Accuracy (overall)	0.74		
Prévalence	0.74	0.23	0.02
Sensibilité	0.98	0.06	0
Spécificité	0.05	0.98	1
Balanced Accuracy	0.51	0.52	0.5

L'accuracy de ce modèle semble être d'environ 74%, mais cela peut être trompeur vu que le déséquilibre de classe n'est pas pris en compte.

Le taux de précision équilibré, offre une représentation plus précise de la performance du modèle

## 1. Arbre de décision

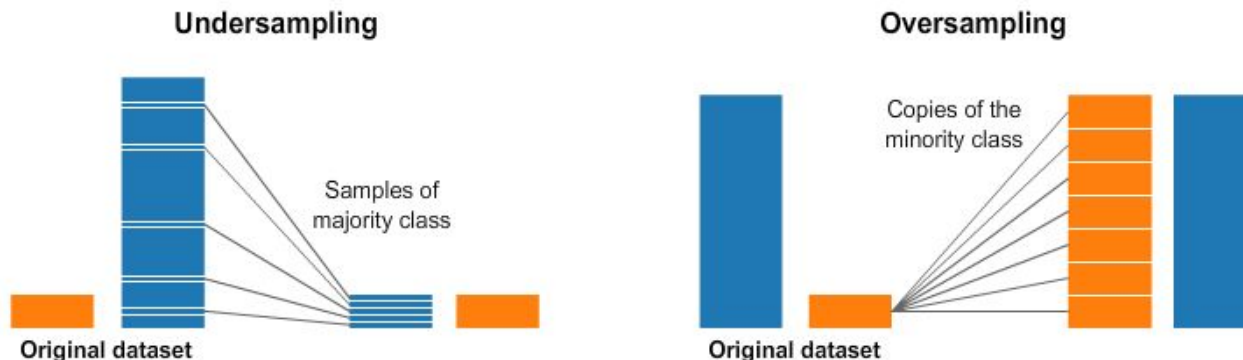
$$\text{Balanced Accuracy} = \frac{\text{Sensitivity} + \text{Specificity}}{2}$$

Where:

- Sensitivity (True Positive Rate) =  $\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$
- Specificity (True Negative Rate) =  $\frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$

## Comment manipuler les jeu de données déséquilibrés?

### 1 - Méthodes d' oversampling & Undersampling



### 2 – Attribuer des poids plus élevés aux échantillons de classes minoritaires.



## 2. Random Forest

**Modèle:** Random Forest (avec des classes pondérées)

**Bibliothèque utilisée:** randomForest

```
rf_model <- randomForest(factor(Music.effects) ~ OCD + Anxiety + Depression + Insomnia,  
  data = train_data,  
  classwt = list(Improve = 1, `No effect` = 2, Worsen = 10))
```

**Matrice de confusion:**

predicted_music_effect	Improve	No effect	Worsen
Improve	94	28	2
No effect	13	5	1
Worsen	0	0	0

**Métrique d'évaluation:**

```
"Accuracy: 0.692307692307692"  
"Balanced Accuracy: 0.510611205432937"
```

## 3. Multinomial Logistic Regression

**Modèle:** Multinomial Logistic Regression

La régression logistique multinomial: prédire une variable catégorielle avec plusieurs niveaux.

**Matrice de confusion:**

	predictions		
	Improve	No effect	Worsen
Improve	204	17	0
No effect	51	12	0
Worsen	6	0	0

**Métriques d'évaluation:**

"Accuracy: 0.744827586206897"

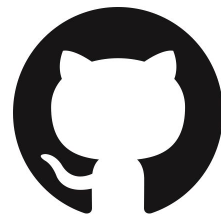
"Balanced Accuracy: 0.371184371184371"

# Portfolio

- En global, il y'a des relations assez logiques dans notre étude.
- Certains facteurs (comme le BPM) ne jouent pas un rôle significatif.
- Jeu de données déséquilibré conduit à des conclusions biaisées du modèle.
- Pistes d'améliorations: plus d'instances pour une analyse plus robuste / ajout de questions sur les données démographiques telles que le pays, le sexe, etc.



# Portfolio



 **DriraYosr** Create README.md 5746a0b · 2 hours ago 🕒 4 Commits

 Music&Health_v1.Rmd	version 1	15 hours ago
 Music&Health_v1.nb.html	version 1 -- version html	15 hours ago
 Music&Health_v2.Rmd	Version 2 du notebook	2 hours ago
 README.md	Create README.md	2 hours ago

 **README** 

Musique et Santé Mentale (MxMH) 🧠🗣️

Le jeu de données provient des résultats de l'enquête sur la Musique et la Santé Mentale sur Kaggle, qui rapporte les résultats sur la préférence des différents genres musicaux et les conditions de santé mentale auto-déclarées (anxiété, dépression, insomnie et trouble obsessionnel-compulsif).



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# Conclusion

- En global, il y'a des relations assez logiques dans notre étude.
- Certains facteurs (comme le BPM) ne jouent pas un rôle significatif.
- Jeu de données déséquilibré conduit à des conclusions biaisées du modèle.
- Pistes d'améliorations: plus d'instances pour une analyse plus robuste / ajout de questions sur les données démographiques telles que le pays, le sexe, etc.



# Merci pour votre attention



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom