

# Zoidberg 2.0

Projet Intelligence Artificielle et Big Data



Introduction

# Notre objectif sur ce projet

Notre but est de développer une solution de machine learning fiable et robuste afin d'assister les médecins dans le dépistage de la pneumonie chez les patients.

Introduction

# Le plan de notre présentation

**Partie #01**

Présentation du dataset et exploration des données.

**Partie #02**

Présentation des différents model de machine learning explorés.

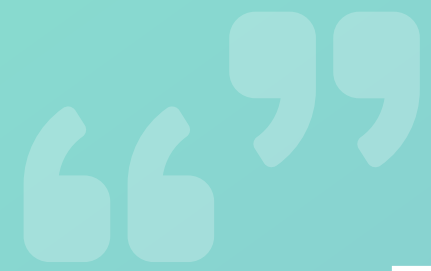
**Partie #03**

Analyse des résultats et comparaison des différentes méthodes.

The background is a gradient from teal on the left to blue on the right. It is decorated with several plus signs of varying sizes and colors (white and light teal).

# Exploration des données

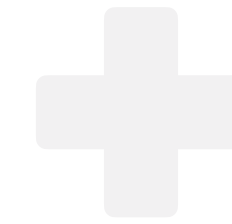
Dataset et traitement des données

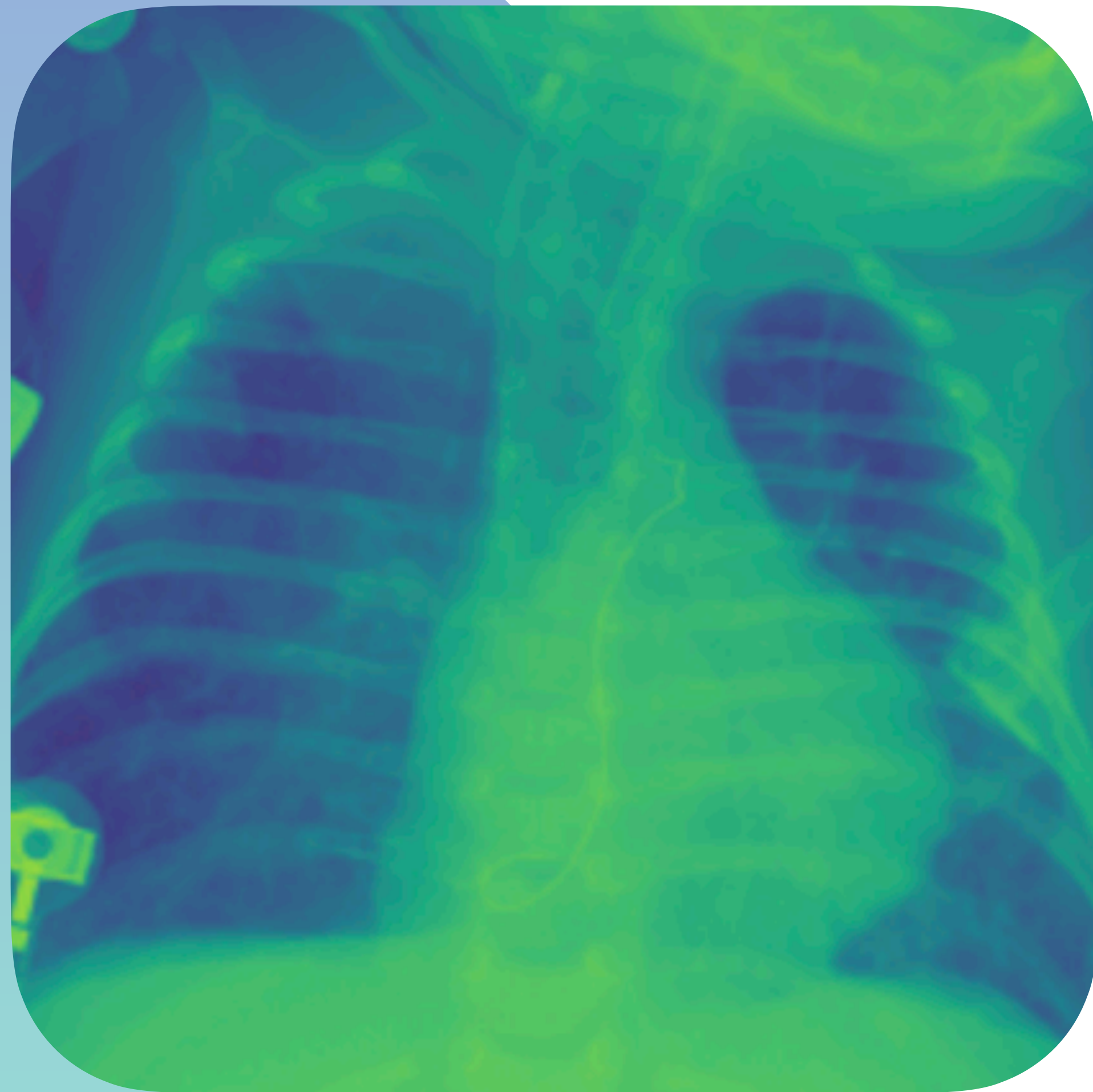


# Dataset

*Un ensemble de données classées par catégories selon certaines caractéristiques prédéfinies.*

*En machine learning, cela permet d'entraîner les modèles.*

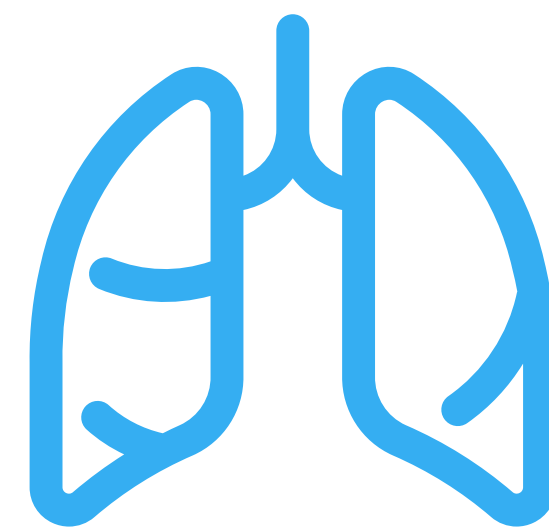




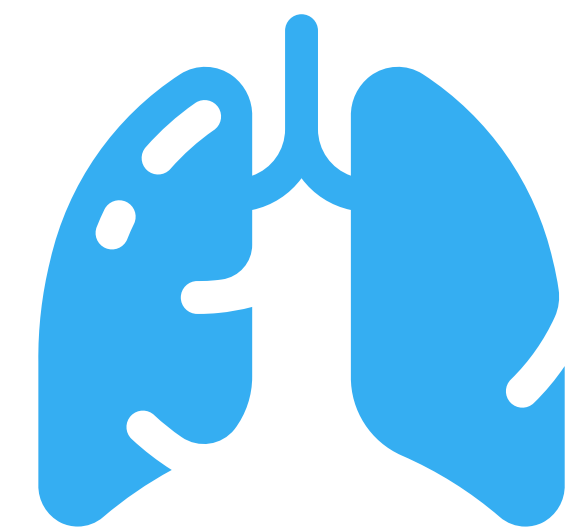
## EXPLORATION DES DONNÉES

# Présentation

Chaque dataset représente un ensemble d'images de radio des poumons:



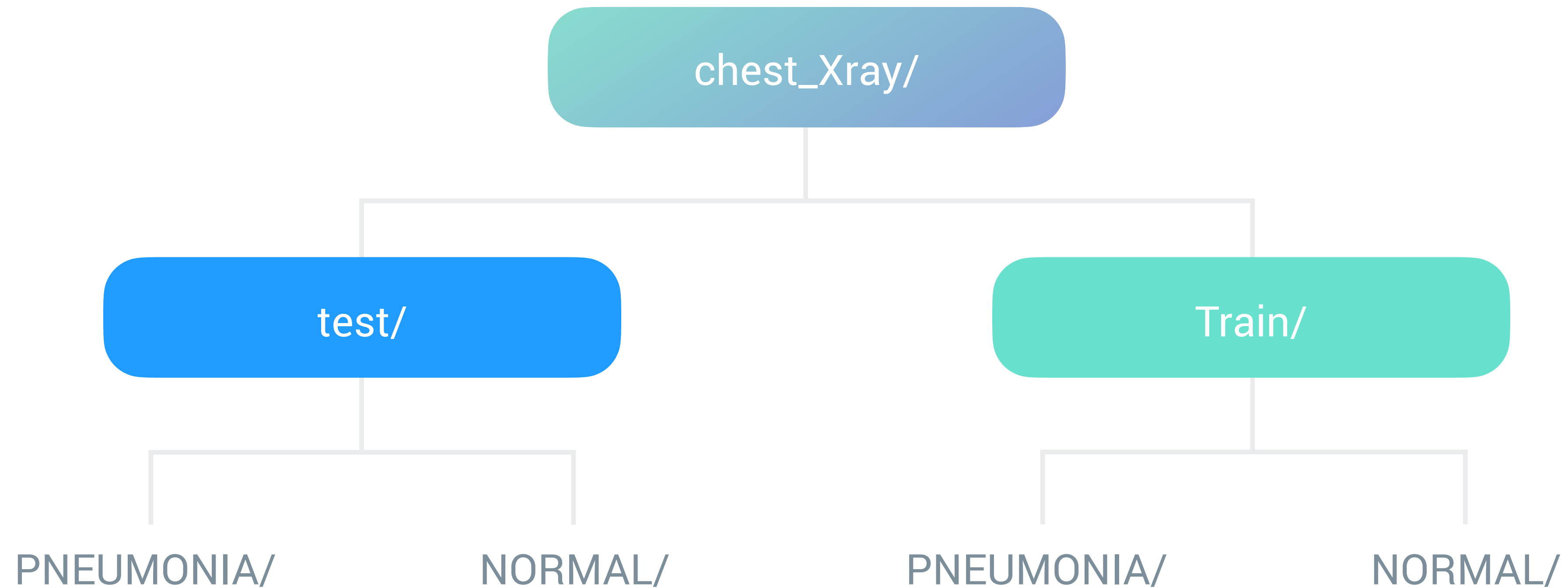
Non infectés (sains)



Infectés (virus ou bactérie)

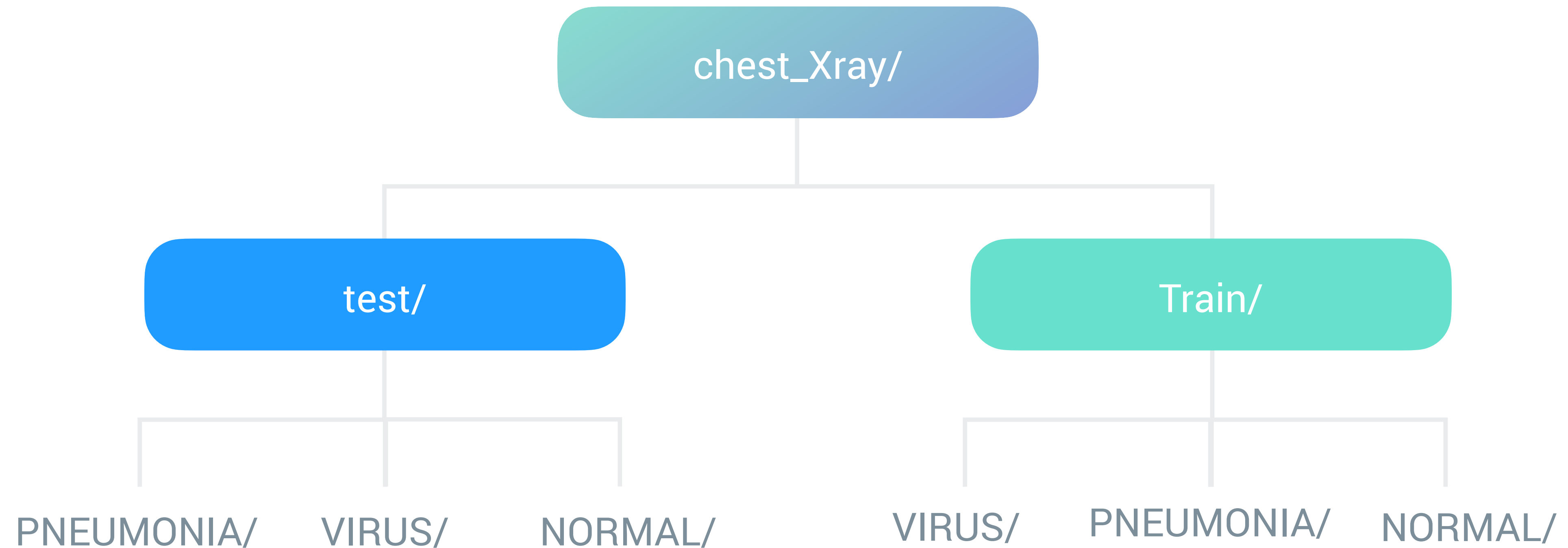
EXPLORATION DES DONNÉES

# Structure du dataset



EXPLORATION DES DONNÉES

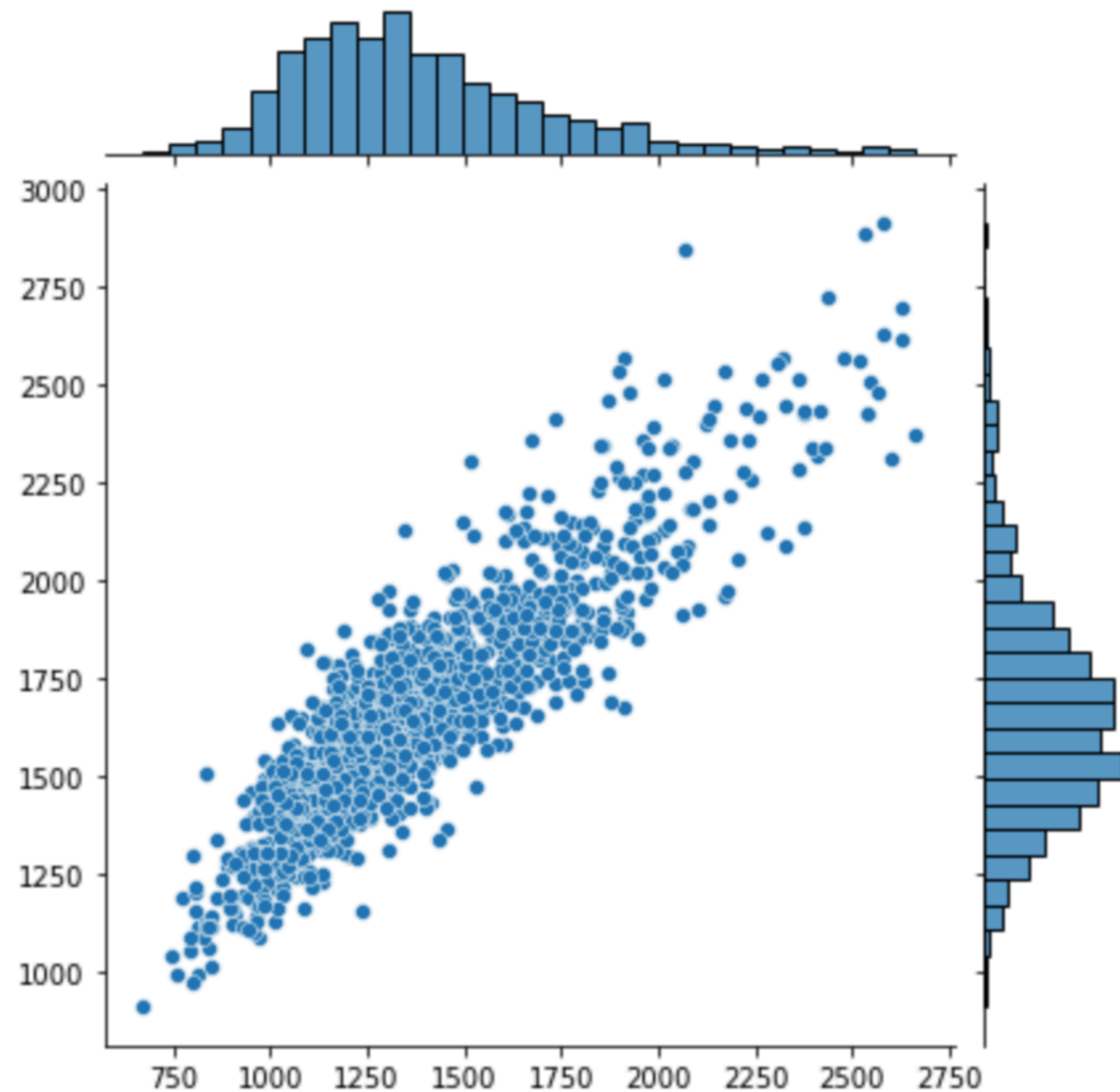
# Nouvelle structure du dataset





## EXPLORATION DES DONNÉES

# Visualisation du format des images



- Trop grande disparité dans les dimensions des images
- Nos models ont besoin d'être entraîné sur des images aux dimensions identiques

EXPLORATION DES DONNÉES

# Traitement des données



## Mise à l'échelle

C'est manipuler les images dans le but d'obtenir les mêmes dimensions



## Normalisation

Faire tenir les valeurs de chaque pixel entre 0 et 1

EXPLORATION DES DONNÉES

# Traitement des données



Diminution des dimensions

Supprimer les pixels qui  
n'apportent pas  
d'informations

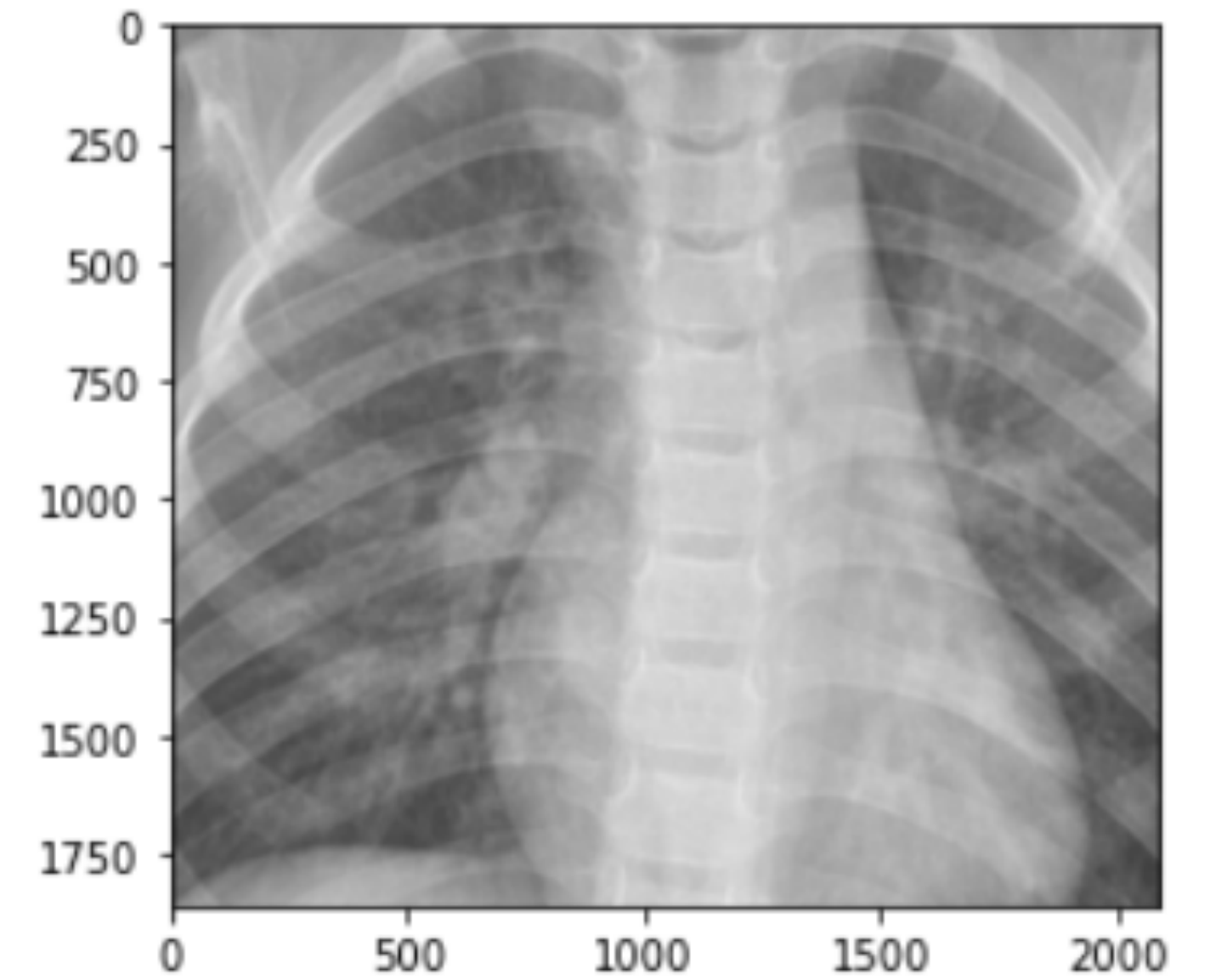
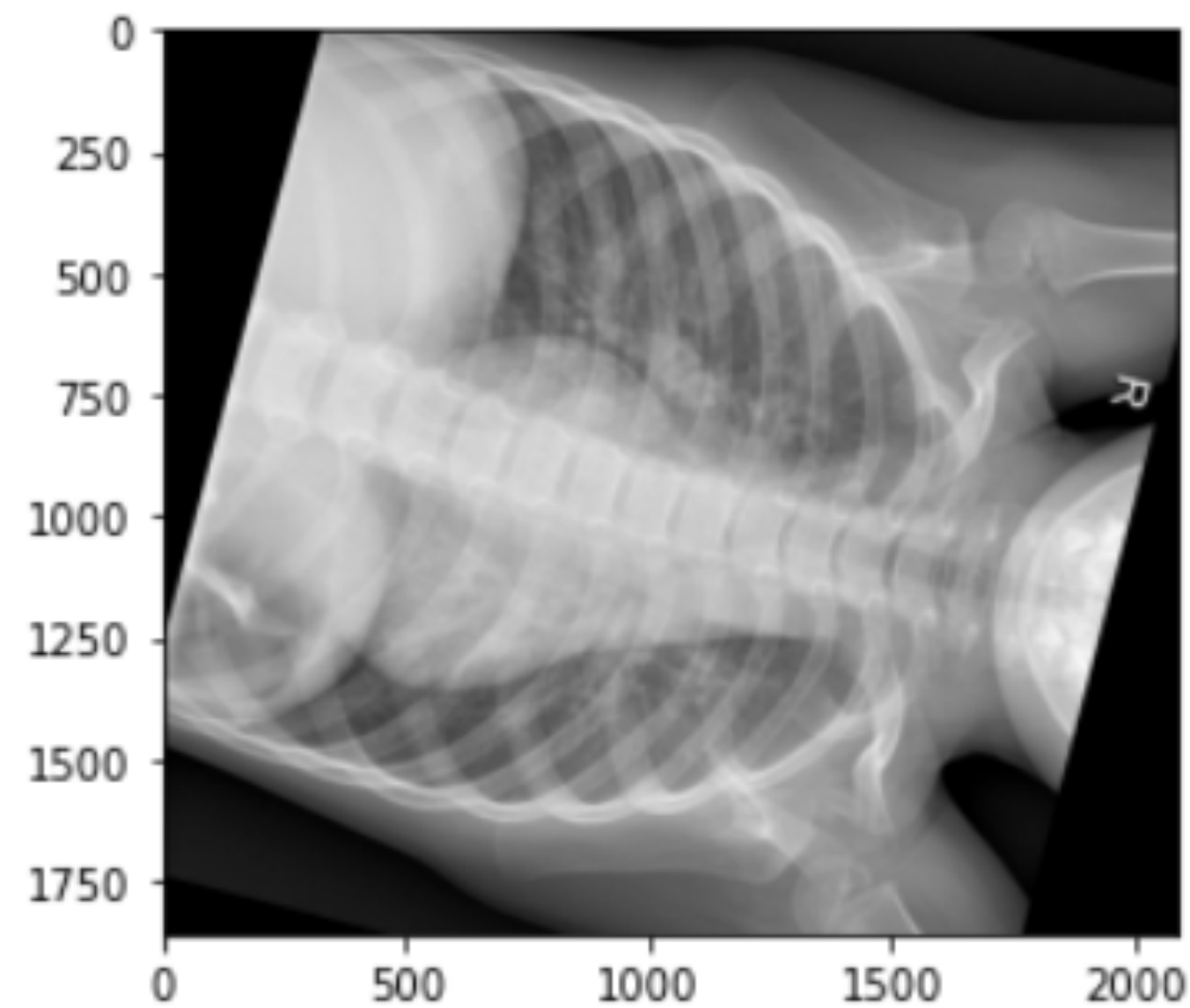
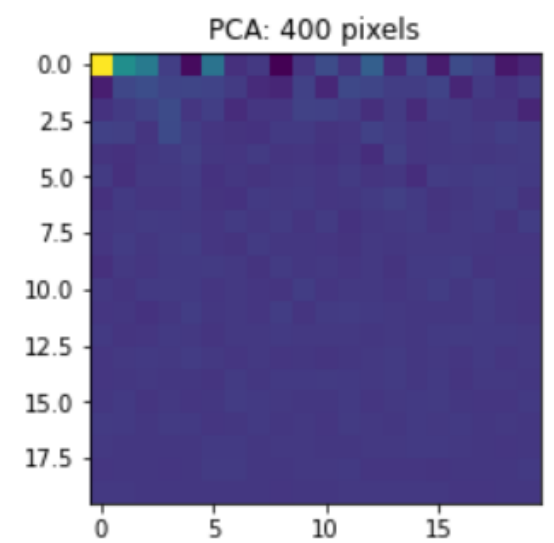
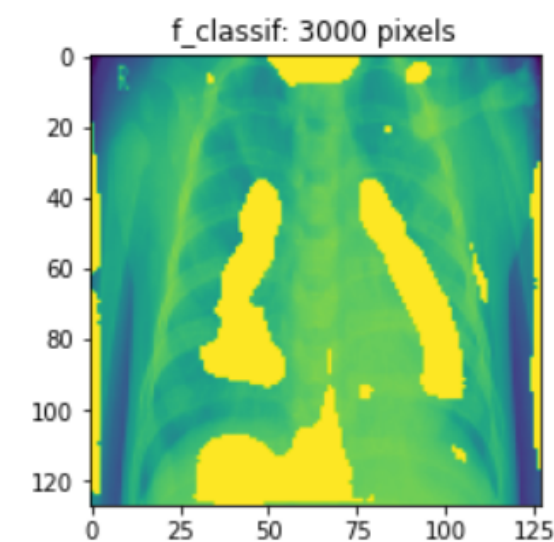
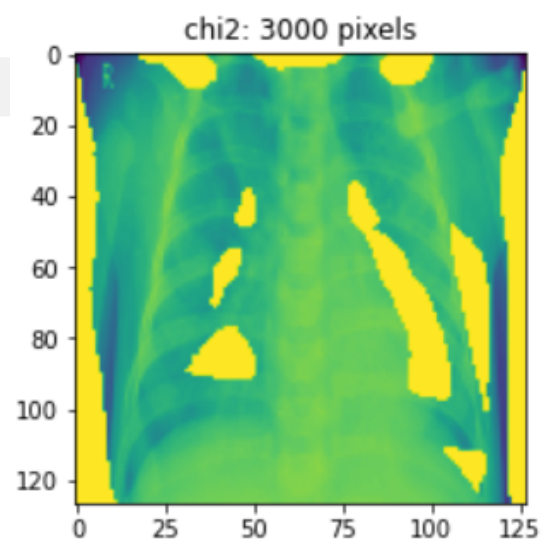
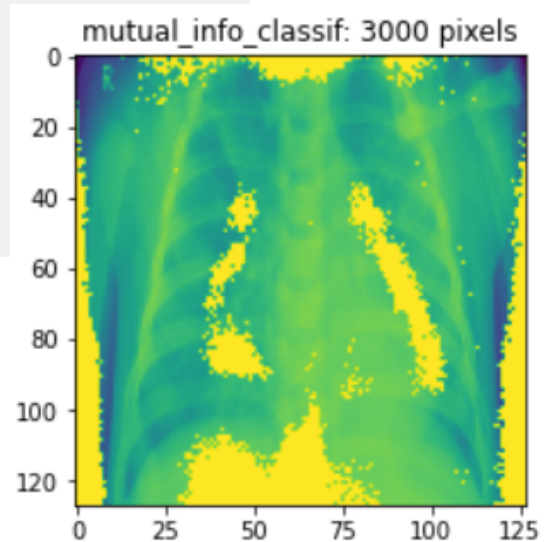


Augmentation des données

Permet de prévenir le sur-  
apprentissage (over-fitting)

EXPLORATION DES DONNÉES

# Traitement des données



The background is a gradient of teal and blue. It features several plus signs of different sizes and colors (light teal, white, and dark teal) scattered across the left side and top.

# Les différents models explorés

Présentation des différentes méthodes de ML utilisées

LES DIFFÉRENTS MODEL EXPLORÉS

# Présentation

Logistic Regression

XGBoost

CATBoost

Neural network

Optimisation  
bayésienne

Transfert de  
connaissances





LES DIFFÉRENTS MODEL EXPLORÉS

# Logistic Regression

- Choisi pour sa rapidité
- Bonne base de prédiction



LES DIFFÉRENTS MODEL EXPLORÉS

# Gradient boosting

- Bon compromis entre les solutions basiques et celles plus gourmandes en ressources
- Facilement interprétable





LES DIFFÉRENTS MODEL EXPLORÉS

# Neural Networks

- Plus complexe à implémenter
- Plus consommateurs en ressources machine
- Plus précis



LES DIFFÉRENTS MODEL EXPLORÉS

# Optimisation Bayésienne

- Trouver automatiquement les meilleurs hyperparamètres
- Economise les ressources



LES DIFFÉRENTS MODEL EXPLORÉS

# Fine tuning

- Utiliser un modèle pré-entraîné
- Gain de ressource et de précision

The background is a gradient of teal and blue. It features several plus signs of different sizes and colors (light teal, white, and dark teal) scattered across the left side and top.

# Analyse des résultats

Présentation des différents résultats obtenus et comparaison des models de ML utilisés



Analyse des résultats

# Metrics

- précision, recall, roc auc, score f1 ?
- le score ROC AUC "macro"

# Tableau de conclusion

Nom du model	ROC AUC	Avantages	Inconniénients	Conclusion
Régression logistique	0.839	Fast	modèle linéaire	Pas pris en compte
XGBClassifier	0.907	Facilement interprétable et précis	Moyennement rapide	Pris en compte
CatBoostClassifier	0.906	Facilement interprétable et précis	Moyennement rapide	Pas pris en compte

# Tableau de conclusion

Nom du model	ROC AUC	Avantages	Inconniénients	Conclusion
Réseau de neurones (SeNet)	0.904	Précis	Pas facilement interprétable/long à entrainer	Pas pris en compte
Réseau de neurones (transmet de connaissance)	0.936	Très précis et rapide à entrainer	Pas facilement interprétable	Pris en compte



# Matrice des métriques et de confusion

Analyse détaillée des résultats



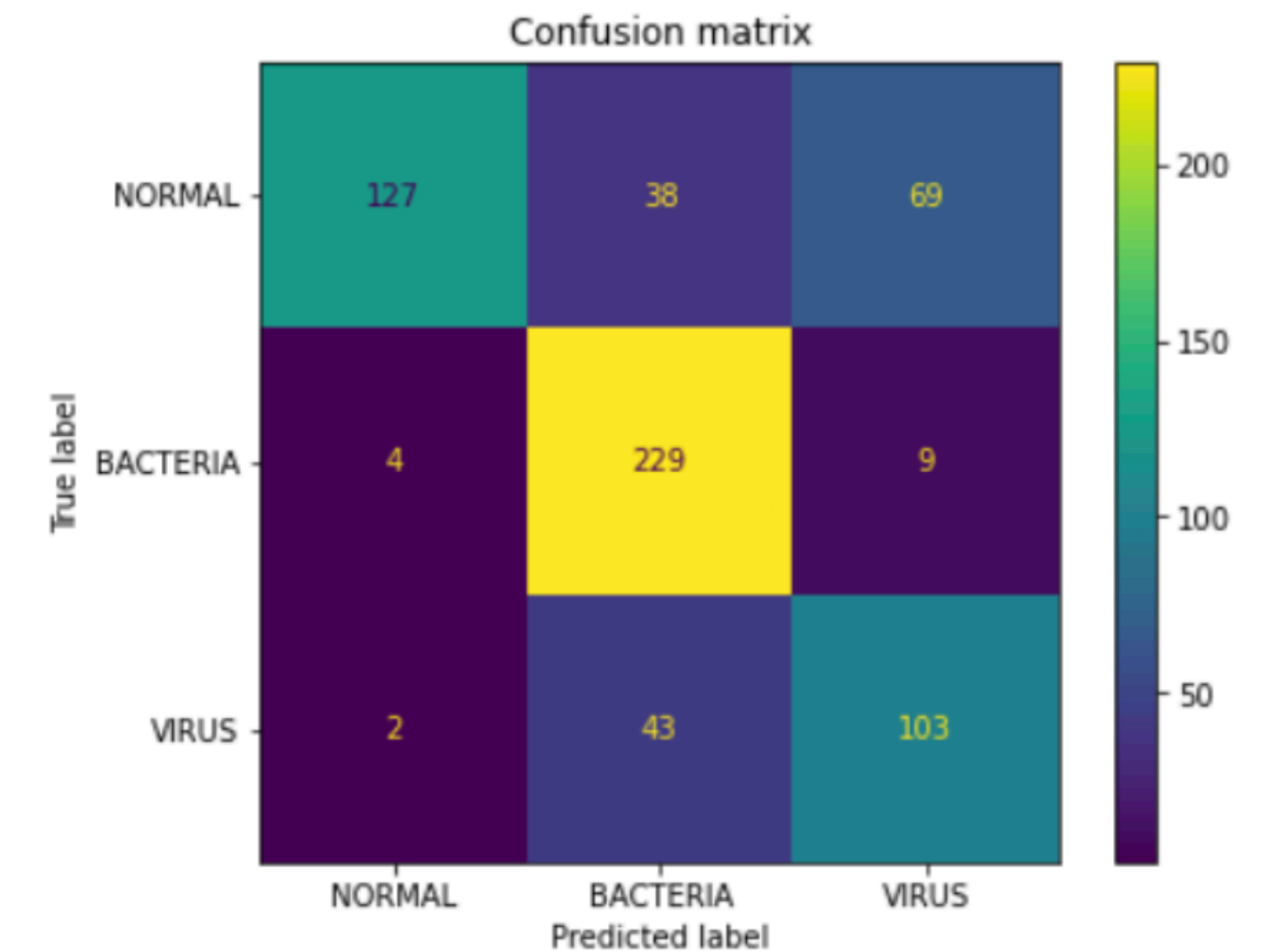
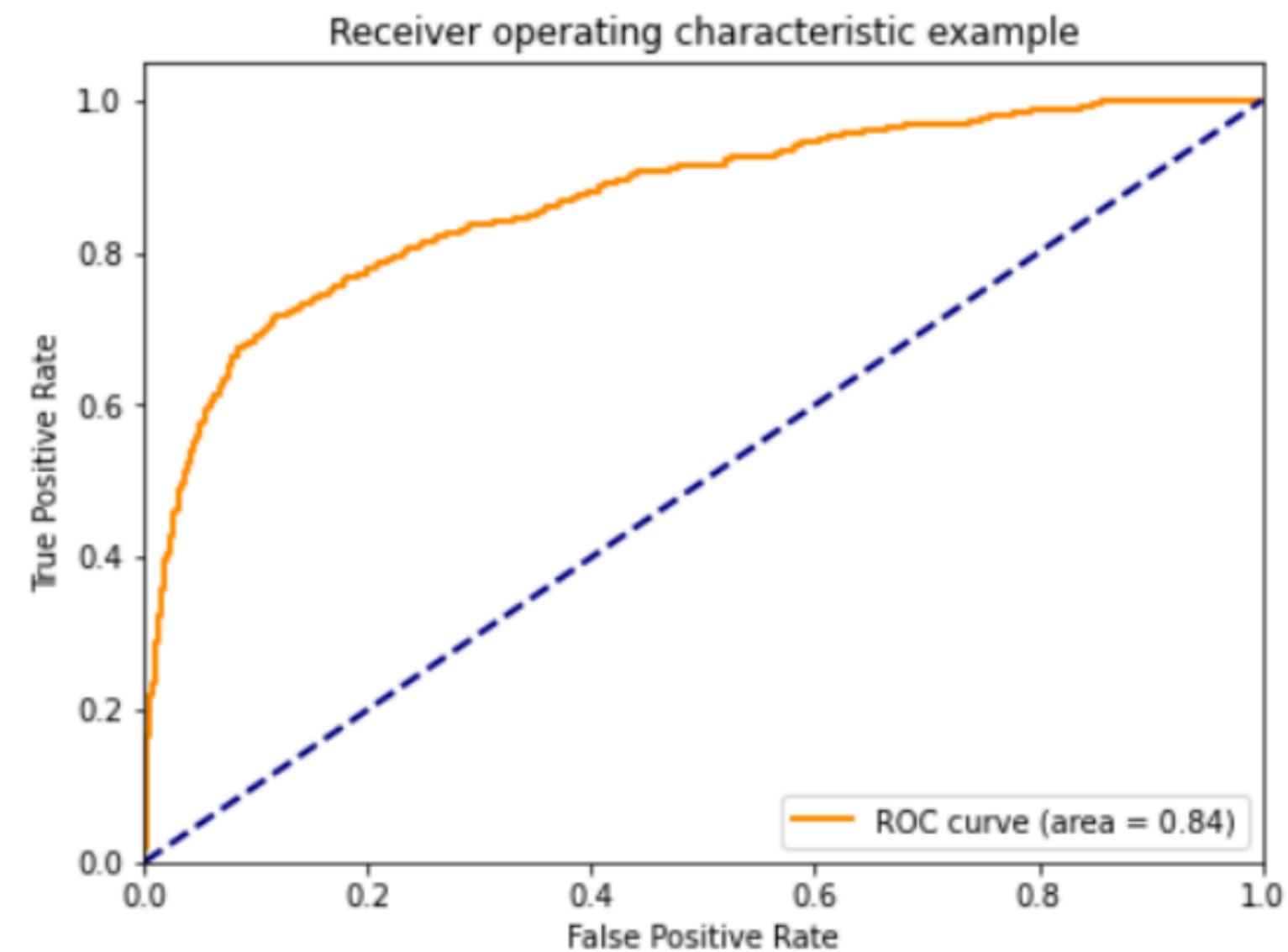
## Matrice des métriques et de confusion

# XGBClassifier

```
classification report:
              precision    recall  f1-score   support

   NORMAL       0.95       0.54       0.69       234
  BACTERIA       0.74       0.95       0.83       242
   VIRUS        0.57       0.70       0.63       148

 accuracy              0.74       624
 macro avg              0.75       624
 weighted avg           0.78       624
```

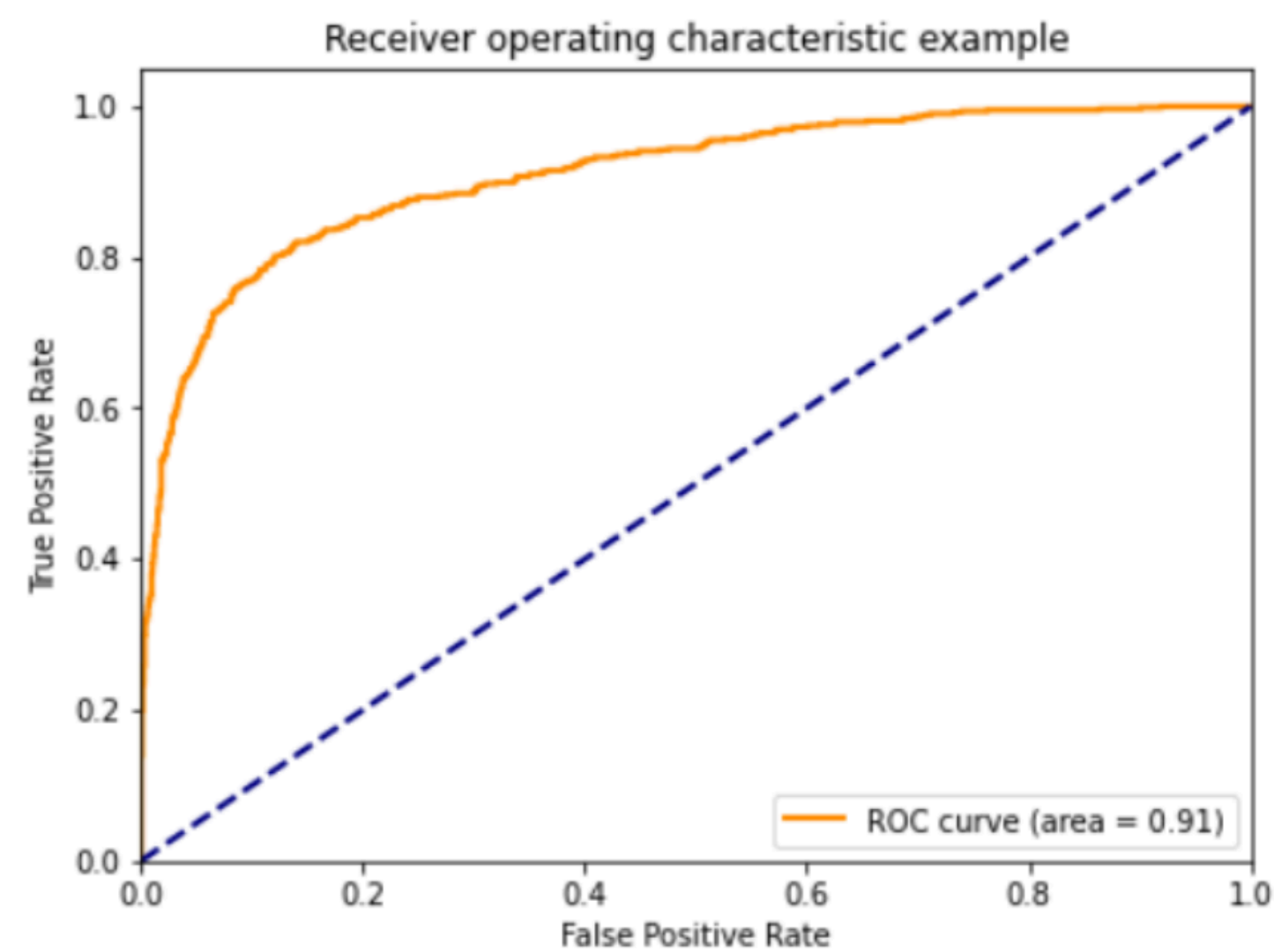


Matrice des métriques et de confusion

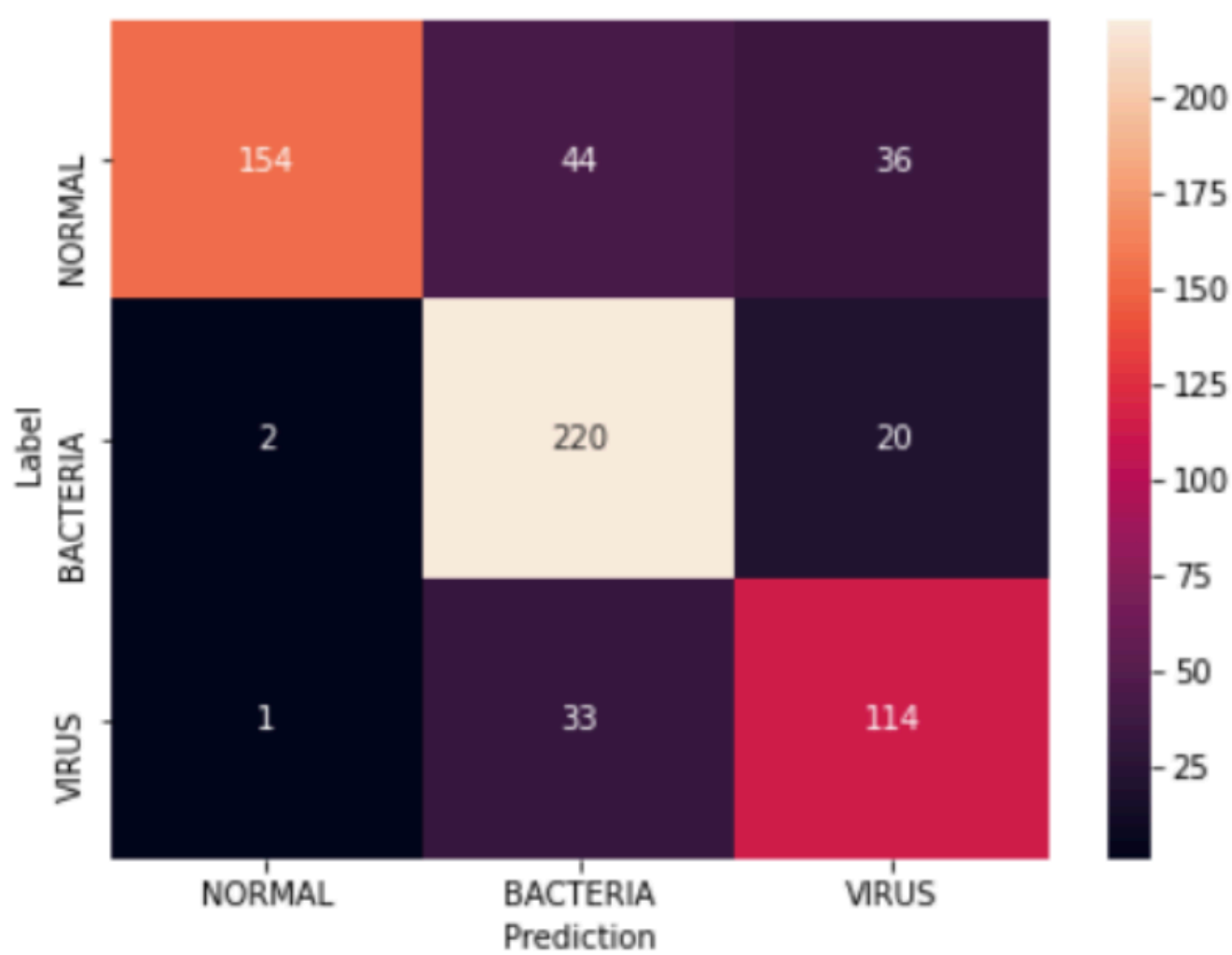
# Réseau de neurone

Voici nos résultats pour notre réseau de neurones:

classification report:				
	precision	recall	f1-score	support
NORMAL	0.98	0.66	0.79	234
BACTERIA	0.74	0.91	0.82	242
VIRUS	0.67	0.77	0.72	148
accuracy			0.78	624
macro avg	0.80	0.78	0.77	624
weighted avg	0.81	0.78	0.78	624



macro ROC AUC: 0.936



The background is a gradient of teal and blue. It features several plus signs of different sizes and colors (white, light teal, and dark teal) scattered across the left side and top. The text is centered on the right side.

# Interprétabilité des résultats

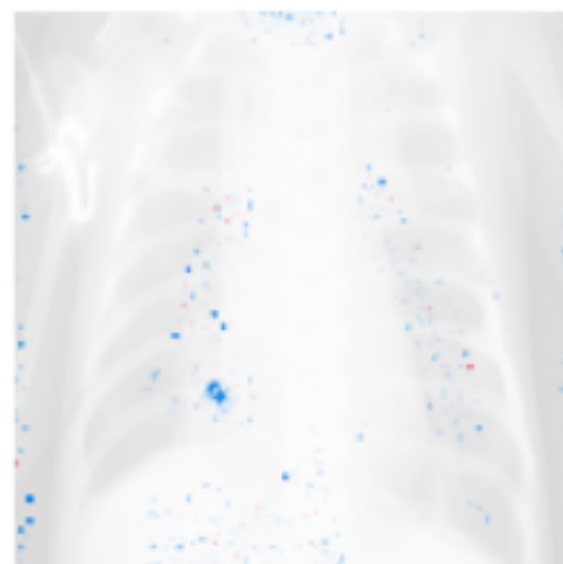
Interprétabilité des résultats

# XGBClassifier

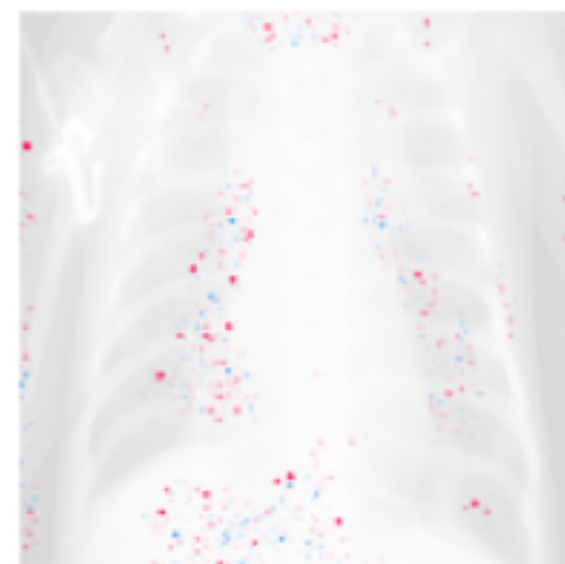
True label: 1



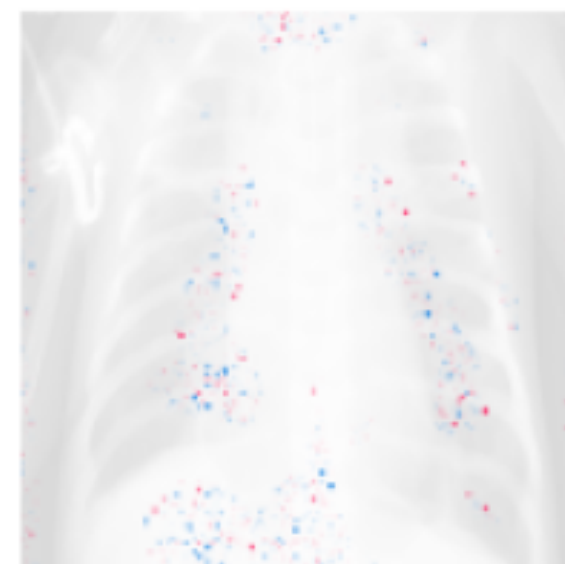
Normal: 0.018%



Bacteria: 0.902%



Virus: 0.080%



-0.015 -0.010 -0.005

SHAP value

0.000 0.005 0.010 0.015

Feature importances



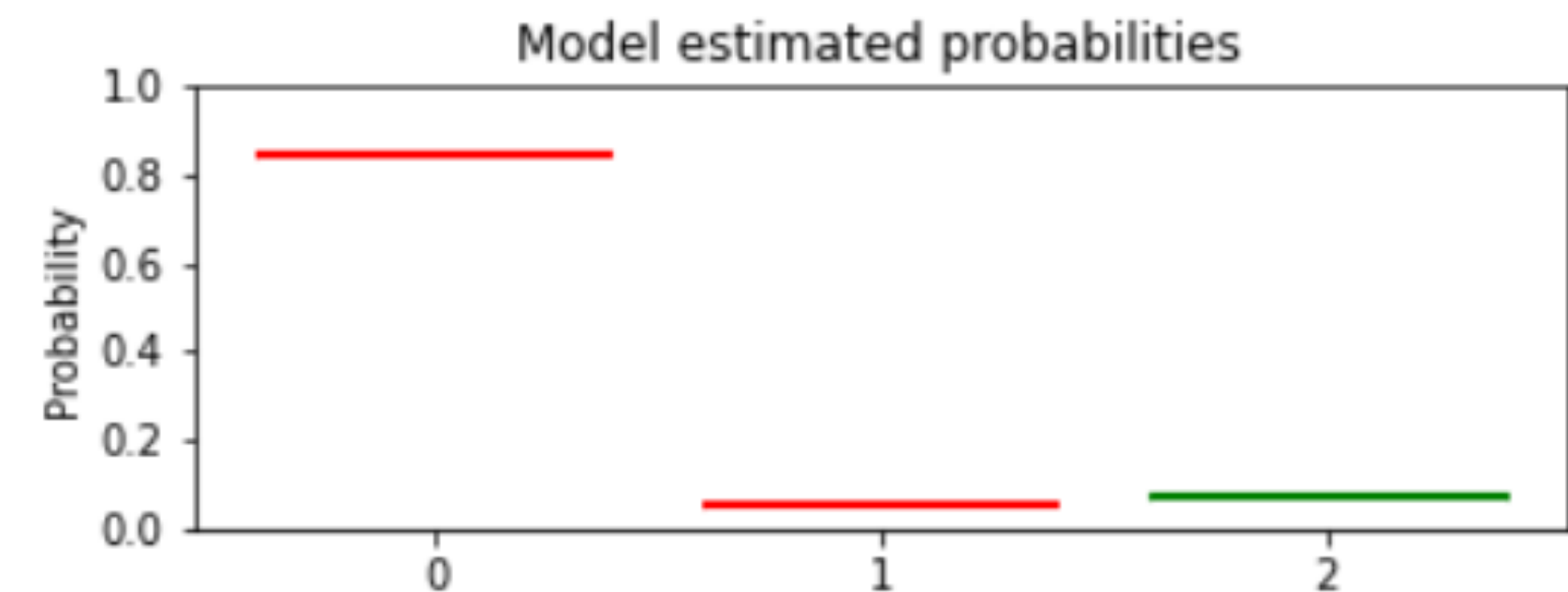
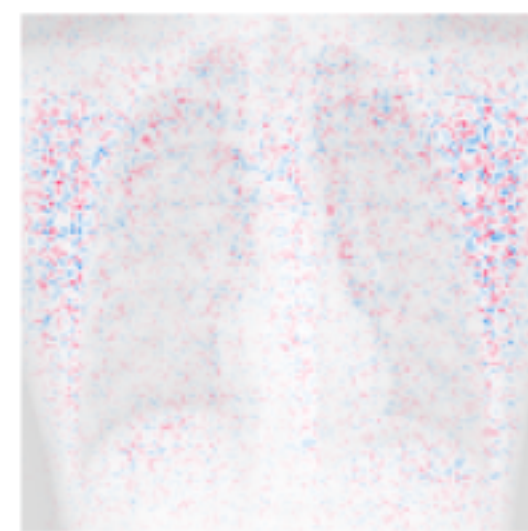
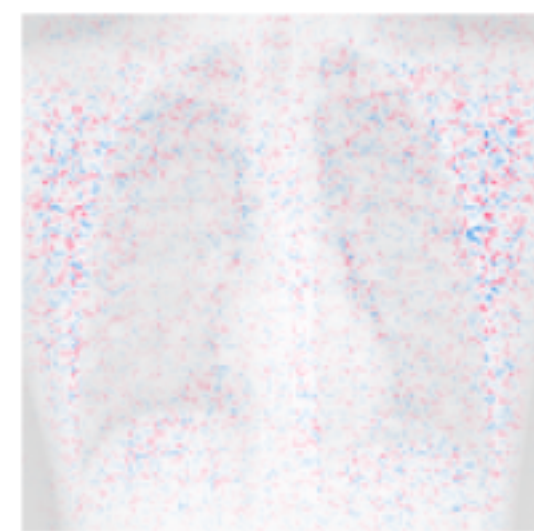
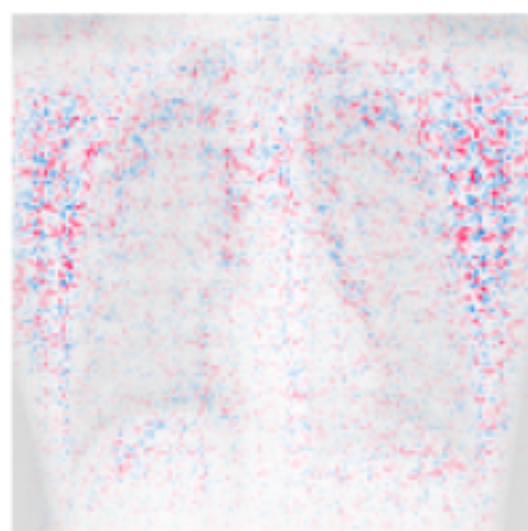
-0.002 -0.001 0.000 0.001 0.002

Feature importances values



Interprétabilité des résultats

# Réseau de neurone



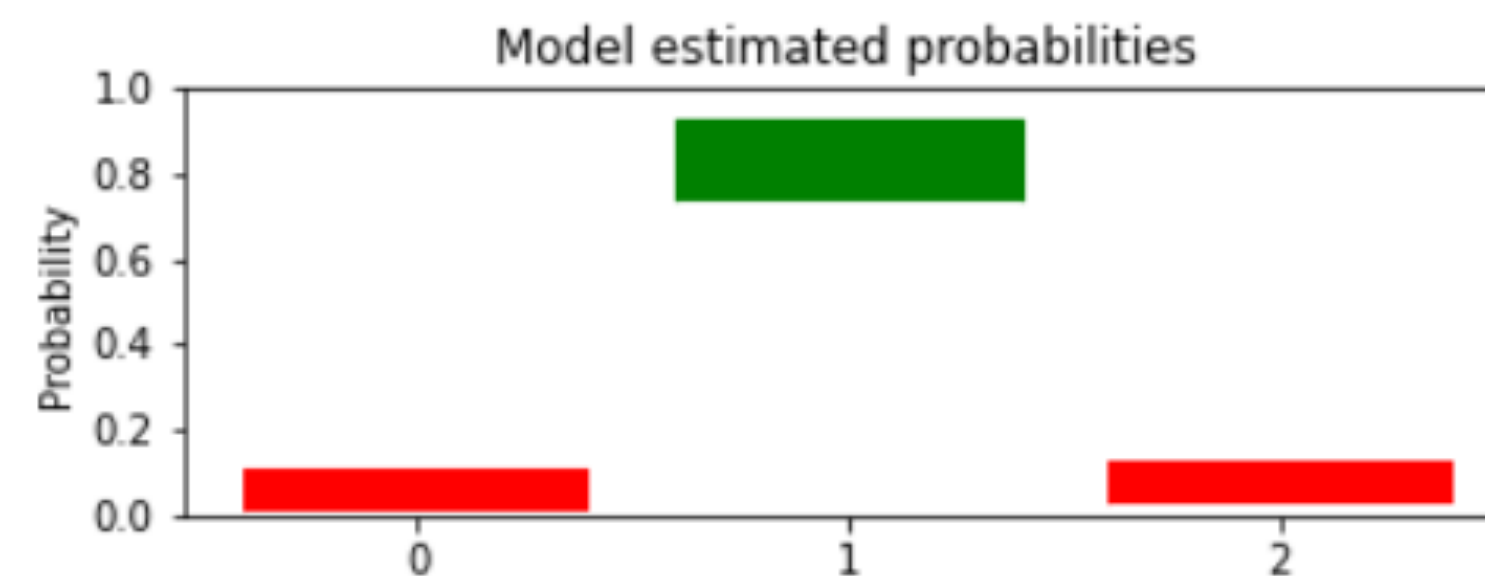
The background is a gradient of teal and blue. It features several plus signs of different sizes and colors (white, light teal, and dark teal) scattered across the left side. The text "Models probabilistes" is centered on the right side in a white, sans-serif font.

# Models probabilistes

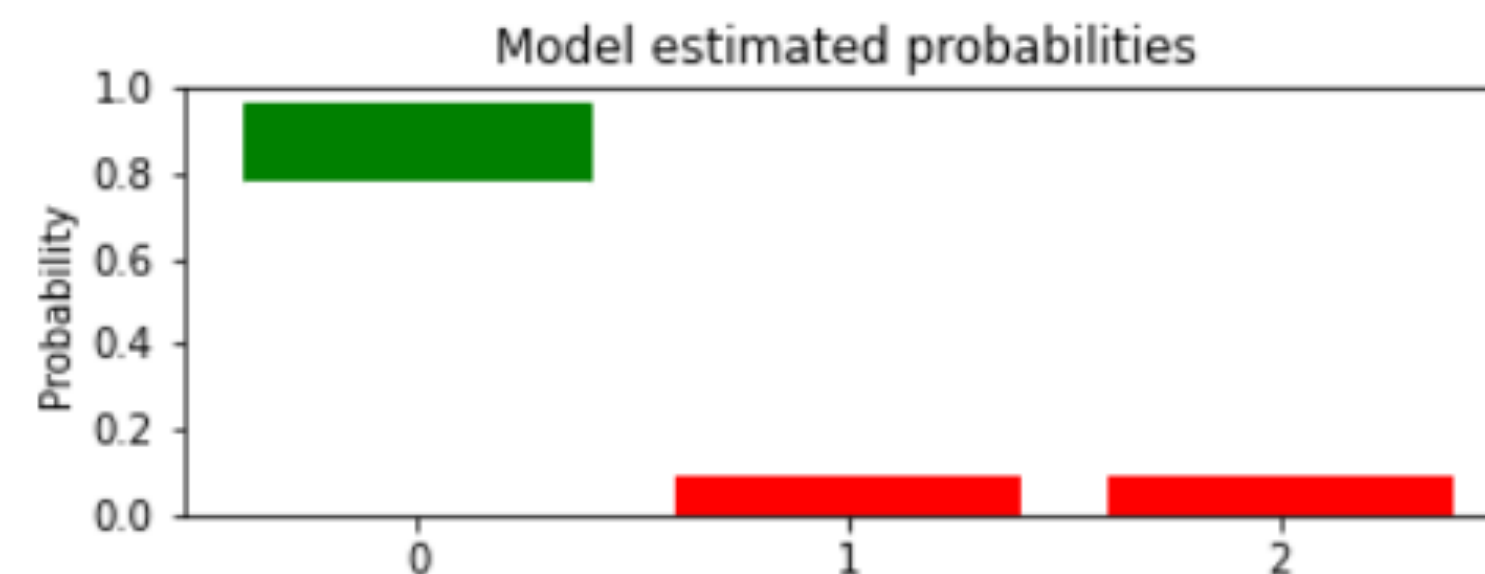
Models probabilistes

# XGBClassifier + bagging

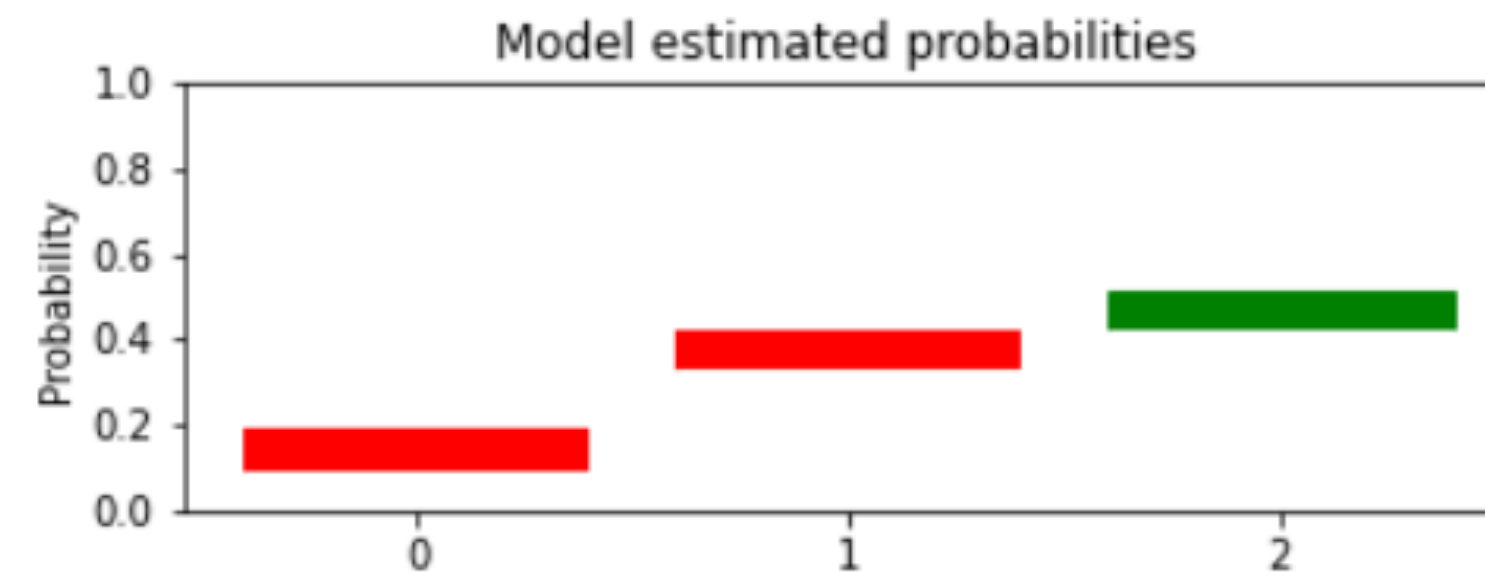
True label: 1



True label: 0



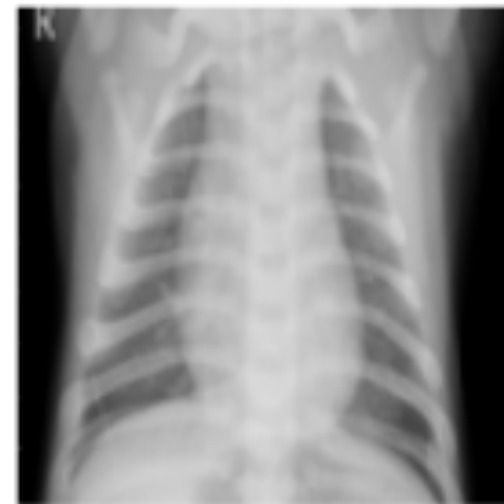
True label: 2



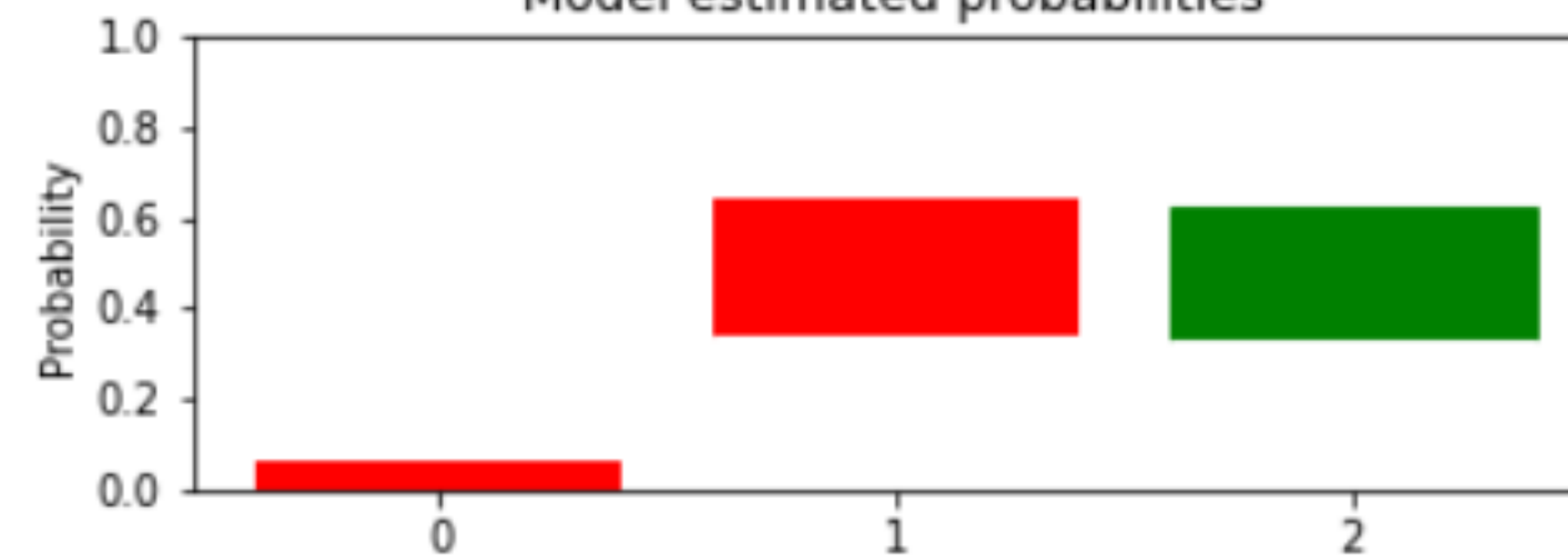
Models probabilistes

# Réseau de neurone

True label: 2



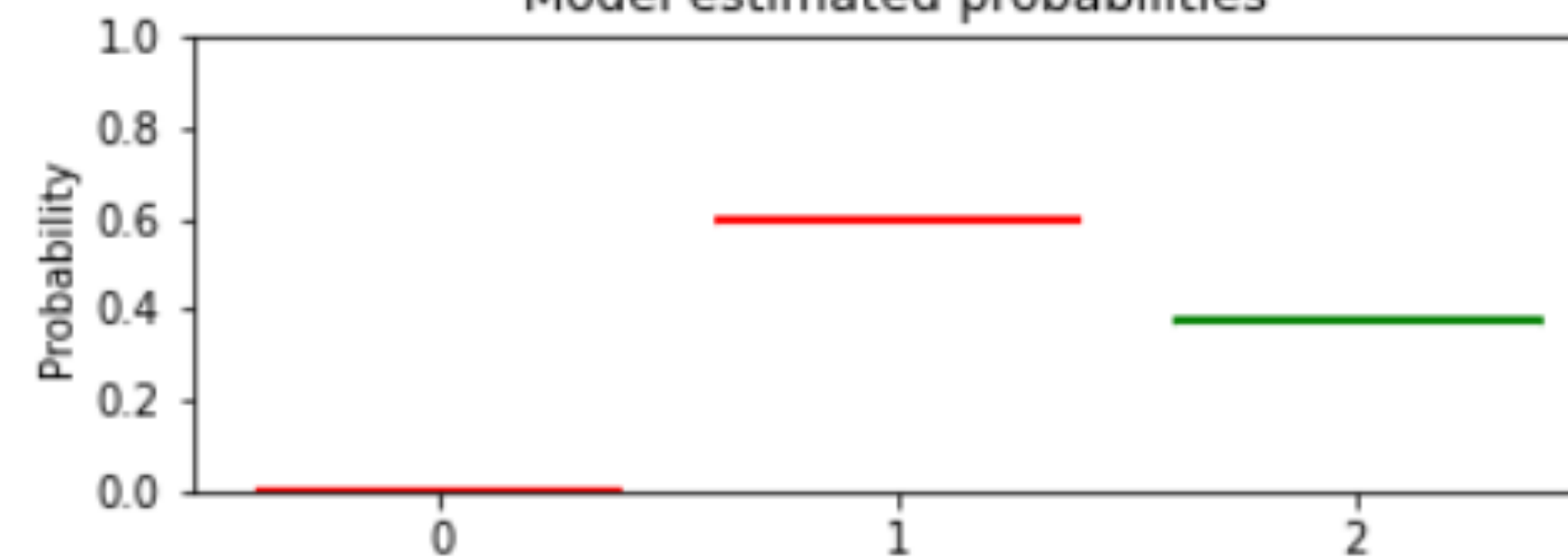
Model estimated probabilities



True label: 2



Model estimated probabilities

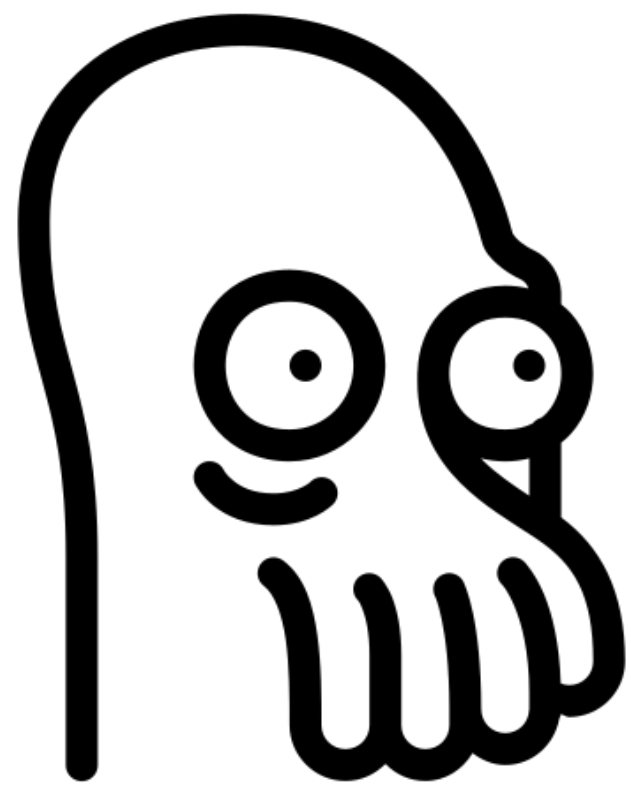




# Conclusion

- Le MobileNetV2 plus performant
- Models classiques plus facilement interprétables





# Merci!

Avez vous des questions?