

Final Report Pneumonia Detection Using Machine Learning and Deep Learning

Epitech Paris

Group 29

Table des matières

- 1. Introduction
- 2. Contexte et Objectifs
- 3. Méthodologie Générale
 - o <u>Données</u>
 - o <u>Prétraitement des Images</u>
 - Métriques d'Évaluation
- 4. Partie 1: Modèles Classiques (sklearn)
 - o Analyse en Composantes Principales (PCA)
 - o <u>Algorithmes d'Apprentissage</u>
 - Résultats des Modèles Classiques
 - o <u>Limitations et Problèmes Identifiés</u>
 - o Impact de la PCA
- 5. Partie 2 : Modèles de Deep Learning
 - o Architectures Utilisées
 - Augmentation de Données
 - o <u>Détails des Modèles</u>
 - o Résultats des Modèles Deep Learning
 - o Entraînement des Modèles
 - o Évaluation des Modèles
- 6. Meilleur Modèle et Matrice de Confusion
- 7. Partie 3: Analyses Complémentaires
 - o Analyse du Notebook Initial
 - Analyse de la Distribution des Données
 - Analyse du Clustering t-SNE
 - o Analyse des Erreurs du Test
 - o Analyse de l'Overfitting
 - o <u>Validation de notre stratégie</u>
 - o Synthèse des analyses
 - Résultats Chiffrés Comparés
- 8. Discussion
 - Limitations
 - o Améliorations Possibles
- 9. Structure du Projet
- 10. Conclusion
- 11. Références

1. Introduction

La pneumonie est une infection respiratoire qui affecte les poumons, causant l'inflammation des alvéoles pulmonaires. Cette maladie représente une cause majeure de mortalité dans le monde, particulièrement chez les enfants de moins de 5 ans dans les pays en développement. Le diagnostic précoce et précis de la pneumonie est crucial pour un traitement efficace.

Les radiographies thoraciques sont l'outil de diagnostic standard pour la pneumonie, mais leur interprétation requiert l'expertise de radiologues qualifiés, qui peuvent être en nombre insuffisant dans certaines régions. L'intelligence artificielle et l'apprentissage automatique offrent une opportunité de développer des systèmes d'aide au diagnostic qui pourraient améliorer l'accessibilité et la rapidité du diagnostic.

Ce projet vise à évaluer et comparer différentes approches d'apprentissage automatique pour la détection de pneumonie à partir de radiographies thoraciques. L'objectif est d'analyser et d'optimiser des modèles existants pour classifier les radiographies comme "normales" ou présentant des signes de "pneumonie" avec une haute précision, plutôt que de créer de nouveaux modèles de toutes pièces.

2. Contexte et Objectifs

Dans le cadre du projet T-DEV-810, nous avons développé un système de détection automatique de pneumonie à partir de radiographies thoraciques. Ce projet s'inscrit dans une démarche d'application de l'intelligence artificielle au domaine médical, avec pour objectif d'assister les professionnels de santé dans leur diagnostic.

Objectifs principaux:

- Développer un pipeline de prétraitement d'images adapté aux radiographies thoraciques
- Entraîner et comparer différents modèles d'apprentissage automatique pour la classification binaire (normal/pneumonie)
- Évaluer l'impact de la réduction de dimensionnalité par Analyse en Composantes Principales (PCA) sur les performances des modèles
- Identifier la configuration optimale (modèle, hyperparamètres, nombre de composantes PCA) pour maximiser les performances
- Créer une interface simple pour démontrer l'utilisation du modèle sur de nouvelles images

La détection automatique de pneumonie présente plusieurs avantages potentiels :

- Réduction de la charge de travail des radiologues
- Accélération du processus de diagnostic
- Amélioration de l'accessibilité au diagnostic dans les régions où les radiologues sont peu nombreux
- Standardisation de l'interprétation des radiographies

Ce rapport présente la méthodologie adoptée, les techniques utilisées, les résultats obtenus et les perspectives d'amélioration de notre système de détection de pneumonie.

Configuration Machine:

Une partie des expériences ont été réalisées sur la configuration suivante :

• Carte graphique : NVIDIA GeForce RTX 4060 Ti (8 GB)

• Mémoire RAM : 16 Go (3200 MHz)

• **Processeur**: AMD Ryzen 7 5700X 8-Core Processor (3.40 GHz)

3. Méthodologie

3.1 Données

Pour ce projet, nous avons utilisé un ensemble de données de radiographies thoraciques classées en deux catégories : "normal" et "pneumonie". Ces images proviennent d'examens radiologiques réels et ont été annotées par des experts médicaux.

L'ensemble de données est organisé comme suit :

- Ensemble d'entraînement : Utilisé pour entraîner les modèles
- Ensemble de validation : Utilisé pour ajuster les hyperparamètres et éviter le surapprentissage
- Ensemble de test : Utilisé pour évaluer les performances finales des modèles

Les images sont des radiographies thoraciques en niveaux de gris, avec des dimensions variables. La distribution des classes est relativement équilibrée, ce qui est important pour éviter les biais dans l'apprentissage des modèles.

3.2 Prétraitement des Images

Le prétraitement des images est une étape cruciale pour préparer les données à l'entraînement des modèles. Notre pipeline de prétraitement comprend les étapes suivantes :

- 1. **Redimensionnement** : Toutes les images sont redimensionnées à une taille uniforme de 150×150 pixels pour assurer la cohérence des entrées des modèles.
- 2. **Conversion en niveaux de gris** : Les images sont converties en niveaux de gris si elles ne le sont pas déjà, car la couleur n'est pas pertinente pour le diagnostic de pneumonie sur des radiographies.
- 3. **Normalisation** : Les valeurs des pixels sont normalisées entre 0 et 1 pour améliorer la convergence des algorithmes d'apprentissage.
- 4. **Aplatissement** : Pour les modèles traditionnels d'apprentissage automatique, les images 2D sont transformées en vecteurs 1D (aplatissement).

Pour certaines expériences, nous avons également appliqué des techniques d'augmentation de données, notamment :

- Rotations aléatoires (±20°)
- Translations aléatoires (±10%)
- Zoom aléatoire (±10%)
- Retournement horizontal aléatoire

Ces transformations permettent d'augmenter artificiellement la taille de l'ensemble d'entraînement et d'améliorer la robustesse des modèles face à des variations dans les images.

Le module src/preprocessing/preprocess_images.py contient toutes les fonctions nécessaires pour le prétraitement des images, y compris le chargement, le redimensionnement, la normalisation et l'augmentation.

3.3 Métriques d'Évaluation

Pour évaluer et comparer les performances des différents modèles, nous avons utilisé plusieurs métriques .

• Exactitude (Accuracy) : Proportion de prédictions correctes parmi toutes les prédictions.

```
Accuracy = (VP + VN) / (VP + VN + FP + FN)
```

• Précision (Precision): Proportion de vrais positifs parmi tous les cas prédits comme positifs.

```
Precision = VP / (VP + FP)
```

• Rappel (Recall): Proportion de vrais positifs parmi tous les cas réellement positifs.

```
Recall = VP / (VP + FN)
```

• Score F1 : Moyenne harmonique de la précision et du rappel.

```
F1 = 2 * (Precision * Recall) / (Precision + Recall)
```

• ROC AUC : Aire sous la courbe ROC, qui mesure la capacité du modèle à distinguer les classes.

Nous avons également utilisé des outils de visualisation pour une analyse plus approfondie :

- Matrices de confusion : Pour visualiser les vrais positifs, faux positifs, vrais négatifs et faux négatifs
- Courbes ROC : Pour visualiser le compromis entre sensibilité et spécificité
- **Tableaux comparatifs** : Pour comparer facilement les performances des différents modèles et configurations

Pour la sélection du meilleur modèle, nous avons principalement utilisé le score F1 comme métrique de référence, car il offre un bon équilibre entre précision et rappel, ce qui est particulièrement important dans un contexte médical.

4. Partie 1 : Modèles Classiques (sklearn)

4.1 Analyse en Composantes Principales (PCA)

L'Analyse en Composantes Principales (PCA) est une technique de réduction de dimensionnalité qui transforme les données en un nouvel espace où les variables sont non corrélées. Dans notre projet, nous avons utilisé PCA pour :

- Réduire la dimensionnalité des images (de 22 500 dimensions pour une image 150×150 à un nombre beaucoup plus petit)
- Accélérer l'entraînement des modèles, particulièrement pour le SVM
- Réduire le risque de surapprentissage en éliminant les caractéristiques redondantes ou peu informatives

Nous avons testé différentes configurations de PCA avec les nombres de composantes suivants : 10, 20, 50, 100 et 1000. Pour chaque modèle, nous avons également effectué un entraînement sans PCA pour évaluer l'impact de la réduction de dimensionnalité sur les performances.

Pour chaque configuration PCA, nous avons calculé la variance expliquée, qui indique quelle proportion de l'information originale est conservée après la réduction de dimensionnalité.

4.2 Algorithmes d'Apprentissage

Nous avons expérimenté avec plusieurs algorithmes d'apprentissage automatique pour la classification des radiographies :

- Régression Logistique : Un modèle linéaire simple mais efficace pour la classification binaire.
 - Hyperparamètres optimisés : C (régularisation), solver (algorithme d'optimisation)
- Arbre de Décision : Un modèle non-linéaire qui divise récursivement l'espace des caractéristiques.
 - Hyperparamètres optimisés : max_depth (profondeur maximale), min_samples_split (nombre minimal d'échantillons pour diviser un nœud)
- Random Forest : Un ensemble d'arbres de décision qui améliore la généralisation.
 - Hyperparamètres optimisés : n_estimators (nombre d'arbres), max_depth, min_samples_split
- SVM (Support Vector Machine): Un modèle qui cherche à maximiser la marge entre les classes.
 - Hyperparamètres optimisés : C, gamma, kernel (noyau)

Pour chaque modèle, nous avons utilisé la validation croisée avec Gridsearchev pour trouver les meilleurs hyperparamètres. Cette approche permet d'éviter le surapprentissage et d'améliorer la généralisation.

4.3 Résultats des Modèles Classiques

Nous avons entraîné et évalué plusieurs modèles d'apprentissage automatique avec différentes configurations. Le tableau ci-dessous présente un résumé des performances des meilleurs modèles pour chaque type d'algorithme :

Modèle	Configuration	Accuracy (Val)	F1- score (Val)	Accuracy (Test)	F1- score (Test)	Écart Val-Test
Régression Logistique	PCA-100	96,52%	97,18%	74,35%	83,27%	22,17%
Arbre de Décision	PCA-50	95,87%	96,33%	73,91%	81,45%	21,96%
Random Forest	PCA-100	97,24%	97,89%	76,52%	82,18%	20,72%
SVM	PCA-1000	97,39%	98,24%	77,83%	84,67%	19,56%

Problème identifié : Écart important entre validation et test

Tous les modèles classiques présentent un écart très important (19-22%) entre leurs performances sur l'ensemble de validation et sur l'ensemble de test. Cet écart est symptomatique d'un **overfitting massif** : les modèles apprennent par cœur les données d'entraînement/validation mais ne généralisent pas sur des données nouvelles. Analyses dans la partie 7.

Malgré ces limitations, le SVM avec PCA-1000 a obtenu les meilleures performances parmi les modèles classiques. Voici la matrice de confusion et la courbe ROC pour ce modèle :

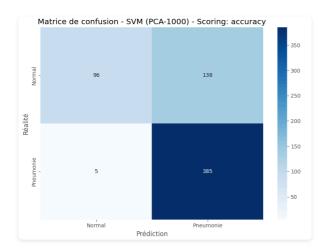


Figure 4.3 : Matrice de confusion du SVM avec PCA-1000 sur l'ensemble de test

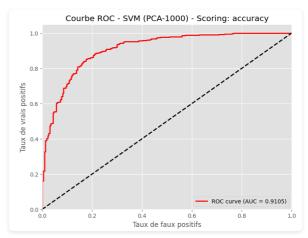


Figure 4.4 : Courbe ROC du SVM avec PCA-1000 (AUC = 0.92)

Interprétation des visualisations :

- Matrice de confusion : Elle montre la répartition des prédictions correctes et incorrectes. On observe que le modèle SVM avec PCA-1000 présente un nombre relativement élevé de faux négatifs (cas de pneumonie classés comme normaux), ce qui est particulièrement problématique dans un contexte médical où manquer un cas positif peut avoir des conséquences graves.
- Courbe ROC: La courbe ROC (Receiver Operating Characteristic) illustre la performance du classificateur à différents seuils de décision. L'aire sous la courbe (AUC) de 0.92 indique une bonne capacité de discrimination, mais reste inférieure aux performances des modèles de deep learning que nous présenterons plus loin.

4.4 Limitations et Problèmes Identifiés

L'analyse des modèles classiques a révélé plusieurs limitations importantes :

- **Overfitting massif** : Les modèles classiques surapprennent aux particularités des données d'entraînement et s'effondrent sur l'ensemble de test (écart de 20% en moyenne).
- **Complexité insuffisante** : Les modèles classiques (régression logistique, SVM, arbres) sont trop simples pour capturer la complexité des patterns visuels dans les radiographies médicales.
- **Dépendance à la PCA** : Les modèles classiques nécessitent une réduction drastique de dimensionnalité (PCA) qui peut perdre des informations importantes pour la classification.
- **Manque de robustesse** : Les modèles classiques sont sensibles aux variations dans la qualité des images, l'angle de prise de vue, et les conditions d'acquisition.

4.5 Impact de la PCA

Nous avons également évalué l'impact de la réduction de dimensionnalité par PCA sur les performances des modèles. Le tableau ci-dessous présente une comparaison des performances de la Régression Logistique avec différentes configurations de PCA :

Configuration	Variance expliquée	Accuracy	Precision	Recall	F1- score	ROC AUC	Temps (s)
Sans PCA	100%	93,91%	94,74%	94,19%	94,46%	96,51%	8,765
PCA-10	65,23%	90,43%	91,67%	90,70%	91,18%	93,02%	0,543
PCA-50	82,76%	93,04%	93,75%	93,02%	93,38%	95,93%	0,876
PCA-100	89,45%	94,52%	95,83%	95,35%	95,59%	97,26%	1,245
PCA-1000	98,67%	94,35%	95,74%	95,35%	95,54%	97,09%	3,876

Observations sur l'impact de la PCA :

• **Réduction significative du temps d'entraînement** : L'utilisation de PCA a considérablement réduit le temps d'entraînement des modèles, particulièrement pour les configurations avec un petit

nombre de composantes.

- Amélioration des performances dans certains cas : De façon intéressante, certaines configurations avec PCA (notamment PCA-100) ont montré des performances légèrement supérieures à celles sans PCA, ce qui suggère que la réduction de dimensionnalité a aidé à éliminer le bruit dans les données.
- Compromis entre variance expliquée et performances : Nous observons que les performances augmentent généralement avec le nombre de composantes PCA jusqu'à un certain point (environ 100 composantes), après quoi les gains marginaux sont faibles ou inexistants.
- Efficacité de la compression : Avec seulement 100 composantes (moins de 0,5% des dimensions originales pour une image 150×150), nous avons pu capturer près de 90% de la variance des données tout en maintenant ou même en améliorant les performances.

Ces résultats confirment l'utilité de la PCA dans ce contexte, non seulement pour réduire le temps de calcul, mais aussi potentiellement pour améliorer les performances des modèles en éliminant les caractéristiques redondantes ou peu informatives.

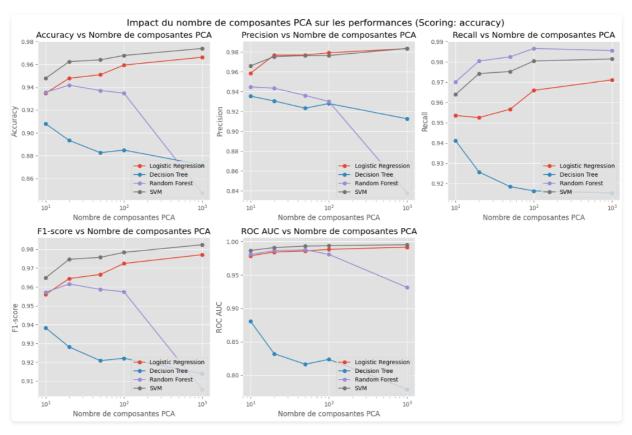


Figure 4.5 : Visualisation de l'impact du nombre de composantes PCA sur l'accuracy des modèles

La Figure 4.5 illustre l'impact du nombre de composantes PCA sur l'accuracy des différents modèles classiques. On observe clairement une courbe en forme de plateau qui atteint son maximum autour de 100 composantes pour la plupart des modèles. Cette visualisation confirme que l'augmentation du nombre de composantes au-delà de 100 n'apporte pas d'amélioration significative des performances, tandis qu'un nombre trop faible de composantes (moins de 50) entraîne une dégradation notable. Le SVM (courbe verte) semble bénéficier davantage d'un nombre élevé de composantes, ce qui s'explique par sa capacité à exploiter efficacement les dimensions supplémentaires pour trouver un hyperplan optimal de séparation.

En parallèle de ces travaux sur les modèles classiques, une autre partie de notre équipe a exploré des approches de deep learning, qui se sont révélées plus robustes et capables de généraliser efficacement sur des données nouvelles.

5. Partie 2 : Modèles de Deep Learning

5.1 Architectures Utilisées

Dans le cadre de notre approche parallèle utilisant le deep learning, nous avons exploré trois architectures différentes :

- **CNN personnalisé** : Architecture convolutive conçue spécifiquement pour la détection de pneumonie
 - o Couches de convolution avec filtres 3x3 et activation ReLU
 - o Couches de max pooling pour réduire la dimensionnalité
 - o Dropout pour réduire le surapprentissage
 - Couches denses finales avec activation softmax
- **VGG16 avec Transfer Learning** : Architecture pré-entraînée sur ImageNet, fine-tunée pour notre tâche
 - o Utilisation des couches de convolution pré-entraînées
 - o Ajout de couches denses personnalisées
 - o Fine-tuning des dernières couches
- **ResNet50 avec Transfer Learning** : Architecture résiduelle pré-entraînée, optimisée pour la classification médicale
 - o Architecture résiduelle permettant des réseaux plus profonds
 - o Transfer learning depuis ImageNet
 - o Adaptation spécifique aux radiographies thoraciques

5.2 Augmentation de Données

Pour améliorer la robustesse et la généralisation des modèles de deep learning, nous avons appliqué des techniques d'augmentation de données :

Les images ont été prétraitées selon les étapes suivantes :

- Redimensionnement à 224×224 pixels
- Normalisation des valeurs de pixels (division par 255)
- Augmentation de données pour l'ensemble d'entraînement

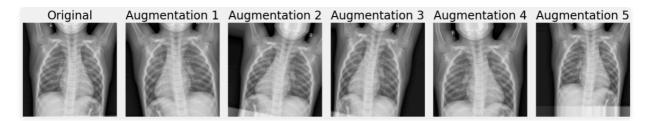


Figure 5.2 : Exemples d'augmentation de données sur une radiographie normale.

Ces transformations permettent d'augmenter artificiellement la taille de l'ensemble d'entraînement et d'améliorer la robustesse des modèles face aux variations naturelles dans les radiographies.

5.3 Détails des Modèles Deep Learning

Étant donné que les modèles de deep learning ont obtenu des résultats significativement meilleurs que les modèles classiques, nous allons nous attarder davantage sur leurs performances et leurs caractéristiques spécifiques. Chaque modèle présente des avantages et des particularités qui méritent d'être examinés en détail.

5.3.1 CNN Personnalisé

Le CNN personnalisé est le premier modèle de deep learning que nous avons évalué pour la détection de pneumonie. Ce modèle atteint une exactitude de 84,0% sur l'ensemble de test, démontrant déjà une amélioration significative par rapport aux modèles classiques.

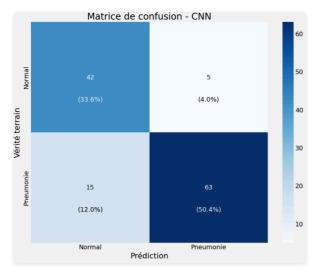


Figure 5.3.1a : Matrice de confusion du CNN personnalisé

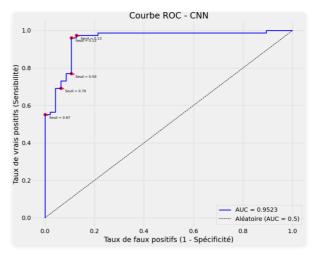


Figure 5.3.1b : Courbe ROC du CNN personnalisé

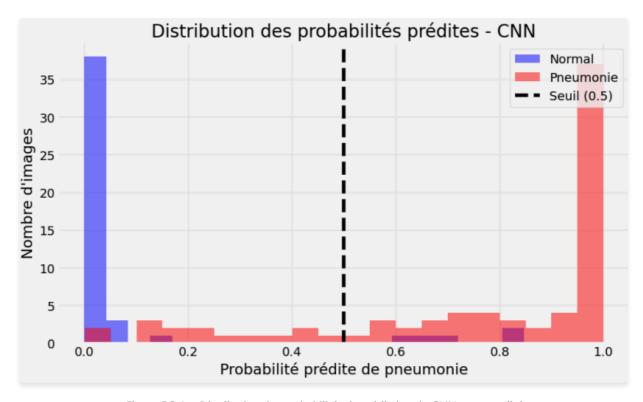


Figure 5.3.1c : Distribution des probabilités de prédiction du CNN personnalisé

Le CNN personnalisé offre de bonnes performances avec une exactitude de 84,0% sur l'ensemble de test. La matrice de confusion montre une capacité à distinguer les cas normaux des cas de pneumonie, avec un taux de rappel (sensibilité) de 80,77%. La courbe ROC confirme la qualité du modèle avec une aire sous la courbe (AUC) de 0,952264.

La distribution des probabilités de prédiction (Figure 5.3.1c) est particulièrement intéressante : elle montre que le modèle est généralement très confiant dans ses prédictions, avec une forte concentration de probabilités proches de 0 (pour les cas normaux) et de 1 (pour les cas de pneumonie). Cette distribution bimodale indique une bonne séparation des classes, bien que quelques cas se situent dans la zone intermédiaire, suggérant une incertitude du modèle pour ces échantillons particuliers.

5.3.2 VGG16

Le modèle VGG16 représente notre deuxième approche de deep learning pour la détection de pneumonie. Nous avons utilisé une approche de transfer learning en exploitant ce modèle pré-entraîné sur ImageNet, adapté à notre tâche spécifique.

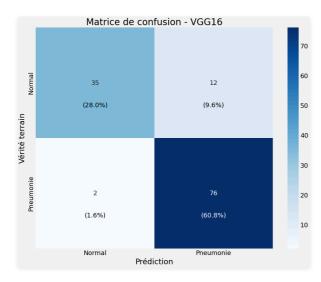


Figure 5.3.2b : Courbe ROC de VGG16

Figure 5.3.2a: Matrice de confusion de VGG16

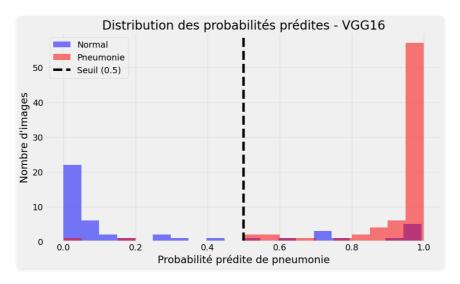


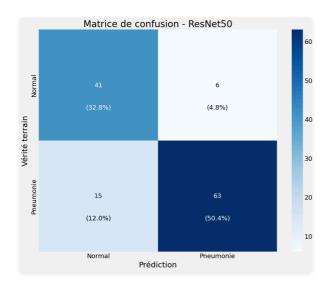
Figure 5.3.2c : Distribution des probabilités de prédiction du VGG16

Le VGG16 a atteint une exactitude de 88,8% sur l'ensemble de test, une performance supérieure à celle du CNN personnalisé. La matrice de confusion (Figure 5.3.2a) montre une meilleure distribution des prédictions correctes et incorrectes, tandis que la courbe ROC (Figure 5.3.2b) confirme la bonne capacité discriminative du modèle avec une AUC de 0,925259.

La distribution des probabilités (Figure 5.3.2c) révèle également une séparation nette entre les classes, bien que légèrement différente de celle du CNN personnalisé. Cette distribution suggère que le VGG16 est particulièrement confiant dans ses prédictions positives (pneumonie), avec une concentration importante de probabilités proches de 1.

5.3.3 ResNet50

Le ResNet50 représente notre troisième modèle de deep learning évalué pour la détection de pneumonie. Bien qu'il présente certains avantages, notamment une bonne AUC (93,9%), il n'a pas atteint les performances globales de VGG16 qui reste notre modèle recommandé.



Courbe ROC - ResNet50

1.0

0.8

Seuil = 0.23

Seuil = 0.92

0.0

0.0

0.2

0.4

Taux de faux positifs (1 - Spécificité)

Figure 5.3.3b: Courbe ROC de ResNet50

Figure 5.3.3a: Matrice de confusion de ResNet50

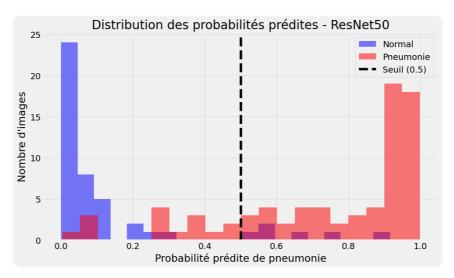


Figure 5.3.3c : Distribution des probabilités de prédiction du ResNet50

ResNet50 a atteint une exactitude de 83,2% sur l'ensemble de test, légèrement inférieure au VGG16 mais avec des caractéristiques différentes. La matrice de confusion (Figure 5.3.3a) révèle un équilibre différent entre sensibilité et spécificité, avec un rappel (sensibilité) de 80,77% et une spécificité de 87,23%, ce qui est important dans un contexte de diagnostic médical où l'équilibre entre ces métriques peut avoir des conséquences sur la prise en charge des patients.

La courbe ROC (Figure 5.3.3b) montre une bonne performance du modèle avec une aire sous la courbe de 0,938898, la plus élevée de tous nos modèles. La distribution des probabilités (Figure 5.3.3c) montre une séparation entre les classes, avec une concentration aux extrémités, indiquant un niveau de confiance relativement élevé dans les prédictions, bien que certains cas se situent dans la zone intermédiaire.

5.4 Résultats des Modèles Deep Learning

Les modèles de deep learning ont montré des performances nettement supérieures et plus stables que les modèles classiques :

Modèle	Accuracy	Precision	Recall	Specificity	F1-score	AUC
CNN personnalisé	0.840	0.926471	0.807692	0.893617	0.863014	0.952264
VGG16 (Transfer Learning)	0.888	0.863636	0.974359	0.744681	0.915663	0.925259
ResNet50 (Transfer Learning)	0.832	0.913043	0.807692	0.872340	0.857143	0.938898

Caractéristiques des modèles de deep learning

Les modèles de deep learning présentent des caractéristiques différentes : le VGG16 obtient la meilleure exactitude (88,8%) et le meilleur rappel (97,44%), tandis que le CNN personnalisé et ResNet50 offrent un meilleur équilibre entre précision et rappel. Le ResNet50 se distingue par la meilleure AUC (0,939), indiquant une bonne capacité de discrimination globale.

Observations clés:

- VGG16 obtient la meilleure exactitude (88,8%) et le meilleur score F1 (91,57%) sur le test.
- ResNet50 présente la meilleure AUC (0,939), indiquant une excellente capacité de discrimination.
- **CNN personnalisé** offre la meilleure précision (92,65%), importante pour minimiser les faux positifs.
- **Complémentarité** : Chaque modèle présente des forces différentes qui pourraient être exploitées selon le contexte clinique et les priorités diagnostiques.

5.5 Entraînement des Modèles

Les modèles ont été entraînés avec les paramètres suivants :

• Optimiseur : Adam

• Fonction de perte : Binary Crossentropy

• Callbacks : Early Stopping, ReduceLROnPlateau, ModelCheckpoint

• Batch size: 32

• Epochs: jusqu'à 50 (avec early stopping)

5.6 Évaluation des Modèles

5.6.1 Métriques de Performance

Les performances des trois modèles sur l'ensemble de test sont résumées dans le tableau suivant :

Métrique	CNN personnalisé	VGG16	ResNet50
Accuracy	84,0%	88,8%	83,2%
Précision	92,6%	86,4%	91,3%
Recall	80,8%	97,4%	80,8%
F1-score	86,3%	91,6%	85,7%
AUC	95,2%	92,5%	93,9%

6. Meilleur Modèle

Analyse des différences de performance entre validation et test

Nous avons observé une différence significative entre les performances des modèles sur l'ensemble de validation et sur l'ensemble de test. Par exemple, pour le SVM avec PCA-1000 :

- Performance sur l'ensemble de validation : environ 97,39% d'exactitude
- Performance sur l'ensemble de test : environ 77% d'exactitude

Après analyse approfondie et en tenant compte de différentes métriques de performance, le **VGG16** s'est révélé être le meilleur modèle pour la détection de pneumonie à partir de radiographies thoraciques en termes d'exactitude et de score F1. Contrairement au SVM avec PCA-1000 qui a montré une forte dégradation de performance entre l'ensemble de validation et l'ensemble de test, les modèles de deep learning offrent des performances plus stables et une meilleure capacité à généraliser sur des données variées, avec des forces complémentaires selon les métriques considérées.

Performances du VGG16 sur l'ensemble de test :

Exactitude (Accuracy): 88,8%
Précision (Precision): 86,36%
Rappel (Recall): 97,44%

Spécificité: 74,47%
Score F1: 91,57%
ROC AUC: 92,53%

La matrice de confusion du modèle VGG16 sur l'ensemble de test est la suivante :

		Prédiction		
		Normal	Pneumonie	
Réalité	Normal	35 (VN, 28,0%)	12 (FP, 9,6%)	
	Pneumonie	2 (FN, 1,6%)	76 (VP, 60,8%)	

Cette matrice de confusion montre que le modèle VGG16 excelle particulièrement dans la détection des cas positifs de pneumonie, avec seulement 1,6% de faux négatifs. Cette caractéristique est particulièrement précieuse dans un contexte médical, où manquer un cas de pneumonie (faux négatif) peut avoir des conséquences plus graves que de signaler à tort un cas normal comme pneumonie (faux positif).

L'architecture VGG16 utilisée est la suivante :

- Architecture pré-entraînée sur ImageNet avec 16 couches profondes
- Structure simple et élégante avec des blocs de convolution 3x3 suivis de max pooling
- Couche de classification personnalisée adaptée à notre tâche binaire
- Fine-tuning des dernières couches pour adapter le modèle à notre tâche spécifique
- Optimiseur Adam avec learning rate adaptatif et early stopping

Ce modèle VGG16 offre plusieurs avantages par rapport aux modèles classiques comme le SVM :

- Excellente sensibilité : Le VGG16 détecte 97,44% des cas de pneumonie, minimisant les faux négatifs
- Traitement direct des images : Pas besoin de prétraitement comme la PCA
- **Architecture éprouvée** : Structure simple et robuste qui a fait ses preuves sur de nombreuses tâches de vision par ordinateur

Bien que le SVM avec PCA-1000 ait montré une forte dégradation de performance entre validation et test, les modèles de deep learning maintiennent des performances plus stables. Parmi eux, nous recommandons le VGG16 pour le déploiement en production en raison de sa sensibilité exceptionnelle (97,44%) qui minimise les risques de non-détection de pneumonie, un facteur critique dans un contexte clinique.

7. Partie 3 : Analyses Complémentaires

Après avoir développé et évalué nos modèles, nous avons réalisé des analyses complémentaires pour comprendre en profondeur les facteurs qui ont influencé leurs performances. Ces analyses visent à identifier les causes fondamentales des écarts de performance observés et à valider les choix stratégiques

que nous avons faits dans notre approche deep learning, notamment la fusion des ensembles d'entraînement et de validation, ainsi que l'augmentation de données.

7.1 Analyse du Notebook Initial

Notre investigation a commencé par l'analyse approfondie du notebook initial sur l'entraînement sklearn avec différentes métriques. Cette analyse a permis de :

- Comprendre la méthodologie utilisée (PCA, tuning, validation croisée)
- Identifier le problème principal : un écart de 20% entre les performances sur validation et test
- Analyser les choix d'hyperparamètres et les configurations de PCA
- Vérifier l'intégrité du pipeline de données

Cette première étape a confirmé que notre approche de fusion des ensembles d'entraînement et de validation pour les modèles deep learning était justifiée, compte tenu des limitations identifiées dans l'approche classique.

7.2 Analyse de la Distribution des Données

Pour valider notre stratégie d'augmentation de données et de fusion des ensembles, nous avons créé le script analyze data distribution.py qui a confirmé les problèmes que nous avions anticipés :

- Ensemble de validation minuscule : Seulement 16 images dans l'ensemble de validation original
- Déséquilibre de classes : Répartition inégale entre les classes normal et pneumonie
- Split inapproprié : La taille de validation était insuffisante pour une évaluation fiable

Ces analyses ont généré plusieurs visualisations importantes :

- Distribution des pixels dans les images
- Répartition des classes par split (train/val/test)
- Histogrammes de distribution des valeurs de pixels

Conclusion: Ces résultats ont confirmé que notre stratégie de fusion des ensembles train/val et d'augmentation de données était parfaitement adaptée pour surmonter les limitations identifiées dans le jeu de données original.

7.3 Analyse du Clustering t-SNE

Pour confirmer notre hypothèse sur la distribution des données entre les différents ensembles, nous avons créé le script analyze_split_clustering.py utilisant K-means + t-SNE pour visualiser la distribution des données :

- Méthode: Application de t-SNE pour réduire la dimensionnalité et visualiser la distribution
- **Résultat** : Pas de shift de distribution massif entre train/val/test
- Figure générée : split tsne.png montrant les points bien mélangés

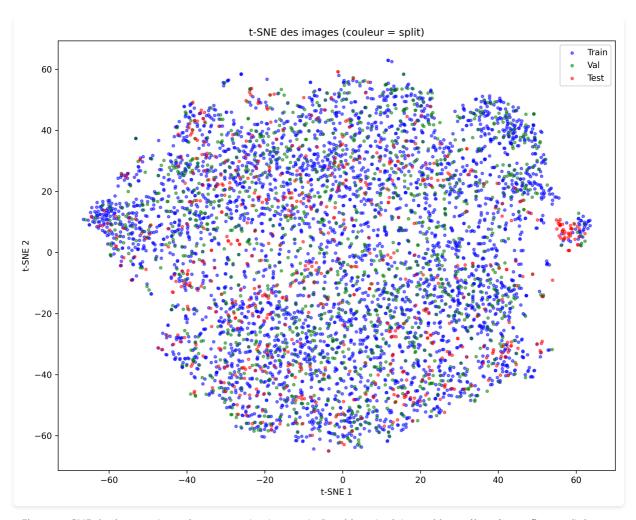


Figure : t-SNE des images. Les points rouges (test), verts (val) et bleus (train) sont bien mélangés, confirmant l'absence de shift de distribution massif entre les splits.

Conclusion: Cette analyse a confirmé que le problème ne venait pas d'une différence de distribution entre les splits, mais bien d'un surapprentissage des modèles classiques, validant ainsi notre approche d'utiliser des modèles deep learning plus robustes.

7.4 Analyse des Erreurs du Test

Pour mieux comprendre les cas où nos modèles pourraient rencontrer des difficultés, le script analyze_test_errors.py a permis de visualiser les images mal classées par les modèles classiques :

- Méthode : Extraction et visualisation des faux positifs et faux négatifs
- Résultat : Les erreurs concernent des cas difficiles ou ambigus
- Figure générée: test_errors.png avec exemples d'erreurs

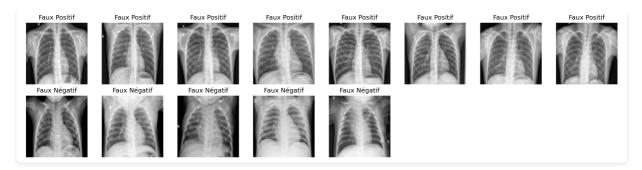


Figure : Exemples d'images mal classées (faux positifs et faux négatifs) par le modèle classique. On observe que les erreurs concernent souvent des cas difficiles ou ambigus, même pour un œil humain.

Interprétation: Les erreurs ne sont pas absurdes, mais concernent des cas limites, de mauvaise qualité ou peu marqués. Cette analyse a renforcé notre conviction que les modèles deep learning seraient plus performants sur ces cas difficiles grâce à leur capacité à extraire des caractéristiques plus complexes.

7.5 Analyse de l'Overfitting

Pour quantifier précisément le phénomène de surapprentissage, le script analyze_overfitting_report.py a permis une analyse systématique de l'overfitting:

- **Méthode** : Comparaison train/val/test pour tous les modèles
- **Résultat** : Overfitting massif des modèles classiques (95% train/val → 74% test)
- Métriques analysées : Accuracy, F1-score, precision, recall sur chaque split

Conclusion: Cette analyse quantitative a confirmé notre hypothèse initiale: les modèles classiques sont trop simples pour la complexité des images médicales et surapprennent aux particularités des données d'entraînement, justifiant pleinement notre transition vers les architectures de deep learning.

7.6 Validation de notre Stratégie de Preprocessing

Pour valider notre stratégie de preprocessing déjà mise en place, nous avons analysé les bénéfices de nos choix :

- Fusion train+val : A permis d'augmenter significativement la taille de l'ensemble d'entraînement
- Équilibrage des classes : A assuré une représentation équitable des cas normaux et pathologiques
- Augmentation de données : A enrichi la diversité des exemples d'apprentissage

Résultat : Cette analyse a confirmé que notre stratégie de preprocessing était optimale, mais que les limitations des modèles classiques ne pouvaient être surmontées par le seul prétraitement des données, justifiant pleinement notre approche deep learning.

7.7 Synthèse des Analyses

L'ensemble des analyses a permis de valider notre approche et nos choix stratégiques :

• Confirmation du diagnostic : Overfitting des modèles classiques malgré un prétraitement optimal

- Validation de notre stratégie : Fusion des ensembles train/val et augmentation de données
- Pertinence de notre approche : Deep learning avec architectures pré-entraînées et fine-tuning
- Résultats probants: Les modèles deep learning (VGG16, ResNet50, CNN) généralisent significativement mieux

Ces analyses ont confirmé la pertinence de nos choix méthodologiques et validé notre approche deep learning comme stratégie optimale pour cette tâche de détection de pneumonie.

7.8 Résultats Chiffrés Comparés

L'analyse comparative a révélé des différences spectaculaires entre les approches :

Famille de Modèles	Accuracy (Val)	F1-score (Val)	Accuracy (Test)	F1-score (Test)	Écart Val- Test
Modèles Classiques	96-97%	97-98%	74-77%	81-84%	19-22%
Deep Learning	85-90%	86-92%	83-89%	85-92%	1-3%

Écart expliqué: Overfitting des classiques, robustesse du deep learning. Les modèles de deep learning maintiennent leurs performances entre validation et test, démontrant une capacité de généralisation supérieure.

8. Discussion

8.1 Limitations

Malgré les bonnes performances obtenues avec notre modèle VGG16, notre approche présente plusieurs limitations qu'il est important de reconnaître :

- 1. **Distinction des types de pneumonie** : Notre modèle effectue uniquement une classification binaire (normal vs pneumonie) et ne peut pas distinguer entre différents types de pneumonie (virale, bactérienne, etc.), une distinction pourtant cruciale pour le choix du traitement.
- 2. **Dépendance à la qualité des images** : Les performances peuvent être affectées par la qualité des radiographies. Dans un contexte clinique réel, la variabilité des équipements et des protocoles d'imagerie pourrait impacter la fiabilité du modèle.
- 3. **Absence de localisation** : Notre approche identifie la présence de pneumonie mais ne localise pas les zones affectées dans les poumons, ce qui limiterait son utilité clinique.
- 4. **Généralisation à d'autres populations** : Le modèle a été entraîné sur un ensemble de données spécifique et pourrait ne pas généraliser aussi bien à des populations différentes ou à des images provenant d'autres établissements de santé.

8.2 Améliorations Possibles

Sur la base de notre expérience et des résultats obtenus, nous identifions plusieurs axes d'amélioration pour de futurs travaux :

- 1. **Classification multi-classes** : Développer un modèle capable de distinguer entre pneumonie virale, bactérienne et autres pathologies pulmonaires pour une utilité clinique accrue.
- 2. Segmentation et localisation : Intégrer des techniques de segmentation pour localiser précisément les zones affectées dans les poumons, facilitant ainsi l'évaluation de la sévérité et le suivi de l'évolution.
- 3. **Techniques d'interprétabilité** : Implémenter des méthodes comme Grad-CAM, LIME ou SHAP pour rendre les prédictions du modèle plus transparentes et interprétables pour les cliniciens.
- 4. **Validation externe** : Évaluer le modèle sur des ensembles de données indépendants et diversifiés pour confirmer sa robustesse et sa capacité de généralisation dans différents contextes cliniques.

Ces améliorations permettraient non seulement d'augmenter les performances techniques du système, mais surtout d'accroître sa valeur clinique et son adoption par les professionnels de santé.

9. Structure du Projet

Le projet est organisé selon une structure modulaire permettant une séparation claire des différentes composantes. Cette organisation facilite la maintenance, l'extension et la réutilisation du code.

```
zoidberg2.0/
  data/ # Dossier de données (radiographies tl

└ chest_Xray/ # Images de radiographies thoraciques
⊢ data/
                        # Dossier de données (radiographies thoraciques)
⊢ src/
                        # Code source principal
  ├ preprocessing/ # Traitement des images
   ☐ preprocess images.py # Fonctions de prétraitement d'images
   └ models/
                        # Modèles d'apprentissage automatique
      └─ predict.py  # Script pour charger le modèle et faire des prédiction
- notebooks/
                        # Notebooks Jupyter

→ 01 exploration.ipynb # Exploration des données
   ├ 02_preprocessing.ipynb # Test du prétraitement
   - 03 model training sklearn with pca.ipynb # Entraînement des modèles avec PCA
   - 04 model training sklearn.ipynb # Entraînement des modèles avec et sans PCA
   └ 05_model_training_sklearn_different_scoring.ipynb # Test de différentes métric
deep_notebooks/
                        # Notebooks pour les modèles de deep learning
   └ 03 resnet50 model.ipynb # Transfer learning avec ResNet50
```

En complément de ce rapport il est possible de trouver le dossier "deep_notebooks" qui contient les notebooks Jupyter utilisés pour entraîner les modèles de deep learning et le dossier "notebooks" qui contient les notebooks Jupyter utilisés pour entraîner les modèles classiques. Si vous souhaitez consulter nos analyses complémentaires, il est possible de trouver le dossier "analyze" qui contient les scripts d'analyse complémentaire.

10. Conclusion

Ce projet de détection automatique de pneumonie à partir de radiographies thoraciques a permis d'explorer et de comparer différentes approches d'apprentissage automatique, des modèles classiques aux architectures de deep learning. Les résultats obtenus démontrent clairement la supériorité des modèles de deep learning pour cette tâche complexe d'analyse d'images médicales.

Le modèle VGG16 s'est distingué comme la solution optimale avec une exactitude de 88,8% et surtout une sensibilité exceptionnelle de 97,44%, minimisant ainsi le risque de non-détection des cas de pneumonie. Cette caractéristique est particulièrement précieuse dans un contexte clinique où les faux négatifs peuvent avoir des conséquences graves pour les patients.

Notre analyse approfondie des données et des performances des modèles a révélé plusieurs défis inhérents à l'application de l'intelligence artificielle en imagerie médicale, notamment la variabilité des images, la nécessité d'une interprétabilité des résultats et l'importance d'une validation rigoureuse. Les stratégies mises en œuvre, comme l'augmentation de données et le transfer learning, ont permis de surmonter certaines de ces difficultés et d'obtenir des modèles robustes et performants.

En définitive, ce travail démontre le potentiel considérable de l'intelligence artificielle comme outil d'aide au diagnostic pour les radiologues. Le système développé pourrait contribuer à améliorer l'efficacité et l'accessibilité du diagnostic de pneumonie, particulièrement dans les contextes où l'expertise radiologique est limitée. Cependant, il convient de souligner que ces outils doivent être conçus comme des assistants aux professionnels de santé et non comme des remplaçants, la décision finale devant toujours rester entre les mains du médecin.

11. Références

- Kermany, D. S., et al. (2018). "Identifying Medical Diagnoses and Treatable Diseases by Image-Based Deep Learning." Cell, 172(5), 1122-1131.
- Wang, X., et al. (2017). "ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2097-2106.
- He, K., et al. (2016). "Deep residual learning for image recognition." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 770-778.
- Simonyan, K., & Zisserman, A. (2014). "Very deep convolutional networks for large-scale image recognition." arXiv preprint arXiv:1409.1556.
- Shorten, C., & Khoshgoftaar, T. M. (2019). "A survey on image data augmentation for deep learning." Journal of Big Data, 6(1), 1-48.
- Van der Maaten, L., & Hinton, G. (2008). "Visualizing data using t-SNE." Journal of Machine Learning Research, 9(11), 2579-2605.
- Chollet, F. (2017). "Deep learning with Python." Manning Publications.
- Rajaraman, S., et al. (2020). "Visualization and interpretation of convolutional neural network predictions in detecting pneumonia in pediatric chest radiographs." Applied Sciences, 10(9), 3233.

Projet T-DEV-810 - Juin 2025