

# AI-Driven Pneumonia Prediction: A Data-Centric Approach Using Machine Learning

Daniel Akama Nyamweya  
SAT 5114

# Problem Statement

- Pneumonia remains a leading cause of morbidity and mortality worldwide, making early detection and prediction is crucial for improving patient outcomes.
- This research explores the application of artificial intelligence (AI) in pneumonia prediction using structured clinical datasets.
- By leveraging machine learning techniques, the study aims to develop a predictive model capable of identifying pneumonia risk based on patient demographics, symptoms, laboratory results, and other relevant features.

## Our Approach: Dimensionality Reduction + Predictive Modeling

- **Goal:** Build a predictive model for pneumonia using patient data.
- **Challenge:** The dataset has 70+ features, which can lead to noise and overfitting.
- **Solution:** Applied PCA to reduce dimensionality while preserving variance.
- **Model:** Training base models i.e. Logistic regression and Ensemble models.

# Code Highlights: Before and after applying PCA to dataset

```
[5]: # Lets print summary information on the dataset
print('Summary information on dataset')
print('-----')
describe_df.data_info()
```

Summary information on dataset

-----

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 768 entries, 0 to 767

Data columns (total 83 columns):

#	Column	Non-Null Count	Dtype
0	record_id	768 non-null	int64
1	age	768 non-null	object
2	gender	767 non-null	float64
3	height	407 non-null	object

Principal Component	Top Influential Features
PC1	comorbid, admission_psi, etio_pneumo_patogen, age
PC2	admission_psi, comorbid, age, sofa_72
PC3	etio_pneumo_patogen, comorbid, discharge_date
PC4	discharge_date, etio_pneumo_patogen, days_ab, weight
PC5	age, weight, height, etio_pneumo_patogen

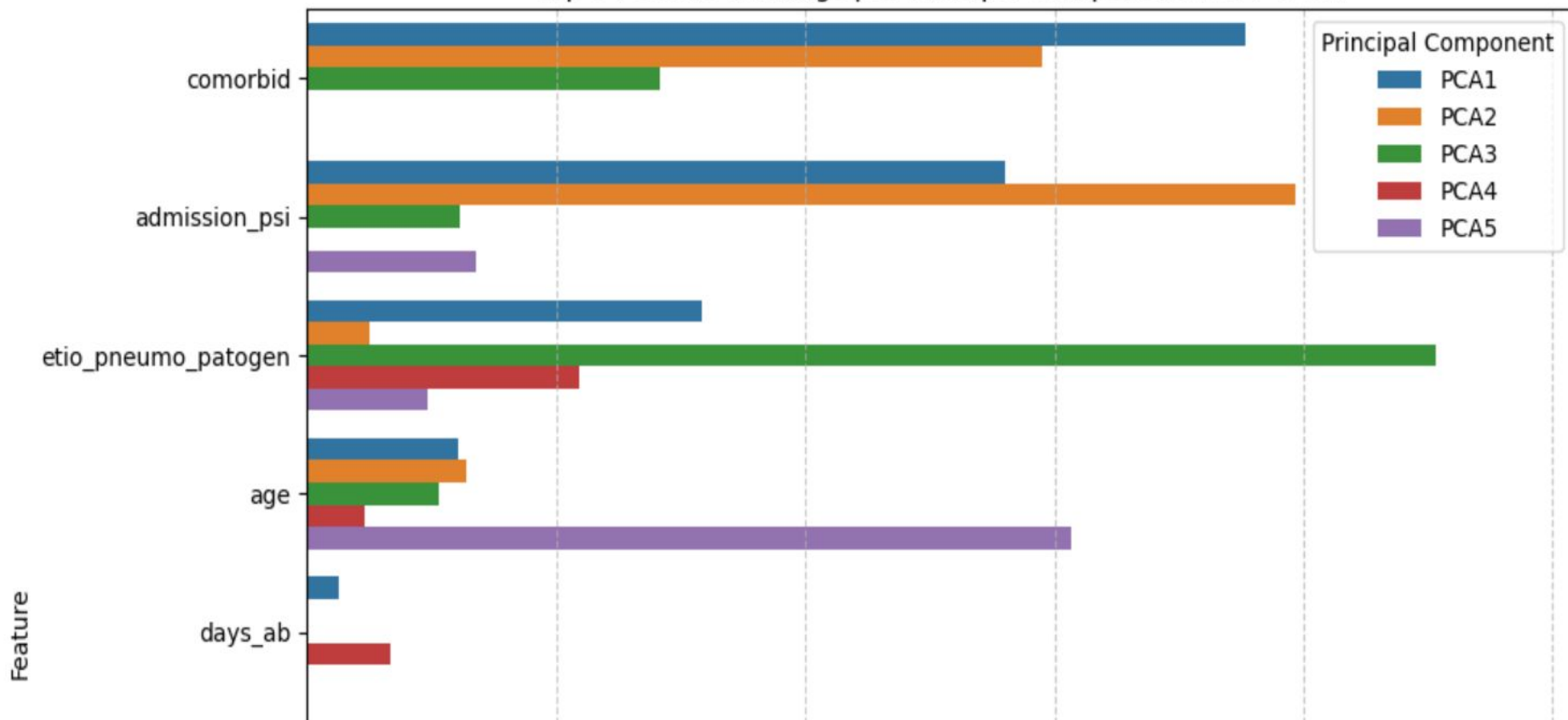
# Performance Metrics

## Model Performance Comparison

Model	Accuracy	Precision (Weighted Avg)	Recall (Weighted Avg)	F1-Score (Weighted Avg)
Logistic Regression	0.933	0.96	0.93	0.95
Random Forest	0.966	0.96	0.97	0.96
XGBoost	0.980	0.96	0.98	0.97

# PCA components

Top 5 Feature Loadings per Principal Component (PC1 - PC5)



# PCA components

Principal Component	Top Influential Features	Interpretation
PC1	comorbid, admission_psi, etio_pneumo_patogen, age	Captures patient condition severity and comorbidity—likely distinguishing patients with complex health profiles. Comorbidity refers to the simultaneous presence of two or more medical conditions in a patient.
PC2	admission_psi, comorbid, age, sofa_72	Emphasizes clinical scores(sofa 72) and patient age—important for assessing initial health status.
PC3	etio_pneumo_patogen, comorbid, discharge_date	Focuses on pneumonia pathogen types and comorbidity—relevant to diagnosis and outcome.
PC4	discharge_date, etio_pneumo_patogen, days_ab, weight	Captures discharge-related info and antibiotic days—indicative of treatment duration or recovery time.
PC5	age, weight, height, etio_pneumo_patogen	Likely represents demographic and anthropometric variability.

# Conclusion

- Models performed well though they significantly favoured the majority class even after applying SMOTE.
- This project was a gateway into advanced AI in Pneumonia research, it proved that clinical data can also be applied to detecting pneumonia to some extent.



# Future Direction

- Combine Structured clinical data with image based data for predictive modeling.
- Apply deep learning techniques such as transfer learning to curb the issue of lack of data.
- Apply Federated Learning, to expose the model to diverse data for enhanced generalizability.
- Deployment. This involves scaling the project to a full-fledged application.

Thank you