

# Predicting\_modeling

May 9, 2025

Myopia Study: Comprehensive Analysis, Modeling and Reporting

---

```
[92]: # Table of Contents

- [0. Previous analysis: Data exploration et initial findings](./
  ↳Data_Exploration.ipynb)
- [1. Introduction: Predictive Modeling and Insights from Previous_
  ↳Analyses](#1-introduction-predictive-modeling-and-insights-from-previous-analyses)
  - [Key Findings and Data Insights from Initial_
  ↳Analyses](#key-findings-and-data-insights-from-initial-analyses)
  - [Model Interpretation Recap](#model-interpretation-recap)
  - [What We Will Address](#what-we-will-address)
- [2. Initialisation](#2-initialisation)
  - [2.1. Importing libraries and loading the myopia dataset for_
  ↳analysis](#21-importing-libraries-and-loading-the-myopia-dataset-for-analysis)
    - [*Data Engineering*](#data-engineering)
- [3. Predicting Modeling](#3-predicting-modeling)
  - [3.1. Feature engineering](#31-feature-engineering)
    - [Commentary](#commentary)
  - [3.2. Selection of variables](#32-selection-of-variables)
    - [3.2.1. VIF Data and Interaction](#321-vif-data-and-interaction)
      - [Commentary](#commentary-1)
    - [3.2.2. PCA](#322-pca)
      - [Commentary](#commentary-2)
    - [3.2.3. Combinations](#323-combinations)
      - [Commentary](#commentary-3)
    - [3.2.3. LASSO and RFE - Automatic Selection_
    ↳Variables](#323-lasso-and-rfe---automatic-selection-variables)
      - [Commentary](#commentary-4)
  - [3.3. Logistic Regression](#33-logistic-regression)
    - [3.3.1 With the variables selected](#331-with-the-variables-selected)
  - [3.4. Results](#34-results)
    - [Commentary](#commentary-5)
    - [3.4.1 Metrics](#341-metrics)
      - [a. Logistic Regression](#a-logistic-regression)
        - [Commentary](#commentary-6)
```

```

- [b. Random Forest](#b-random-forest)
  - [Commentary](#commentary-7)
- [c. GradientBoosting](#c-gradientboosting)
  - [Commentary](#commentary-8)
- [d. xgboost](#d-xgboost)
  - [Commentary](#commentary-9)
- [e. Comparison](#e-comparison)
- [Comparative Analysis of Model Calibration and Precision-Recall
↪Performance](#comparative-analysis-of-model-calibration-and-precision-recall-performance)
- [4. **Global Synthesis and Conclusion**](#4-global-synthesis-and-conclusion)
  - [**Key Achievements et Insights**](#key-achievements--insights)
  - [**Comparative Model Evaluation**](#comparative-model-evaluation)
  - [**Clinical and Practical
↪Implications**](#clinical-and-practical-implications)
  - [**Conclusion**](#conclusion)

```

The history saving thread hit an unexpected error (OperationalError('attempt to write a readonly database')).History will not be written to the database.

```

Cell In[92], line 3
- [0. Previous analysis: Data exploration et initial findings](./
↪Data_Exploration.ipynb)
~
SyntaxError: invalid syntax. Perhaps you forgot a comma?

```

## 1 1. Introduction: Predictive Modeling and Insights from Previous Analyses

In this section, we leverage the results and findings from the initial ‘Predicting\_modeling’ analysis ([notebook link provided in the project files](#)). Our modeling work is grounded in a comprehensive exploratory and statistical investigation (see previous sections), which has helped us identify key risk factors and relationships relevant to myopia.

### 1.0.1 Key Findings and Data Insights from Initial Analyses

**1. Features strongly associated with myopia:**

- **SPHEQ (Spherical equivalent refraction):** The single most discriminative variable for myopia. Statistically significant (very low p-value), consistently top-ranked by feature importance in all models. Lower SPHEQ values are strongly predictive of myopia.
- **PARENTSMY (At least one myopic parent):** Strong and highly significant association with myopia incidence. Parental (hereditary) status drastically increases the risk for childhood myopia. This is supported by proportion tests and feature importance analysis.
- **SPORTHR (Hours in sports/outdoors):** Statistically significant; lower outdoor activity is associated with myopia.

**2. Features not associated or weakly associated with myopia:**

- **GENDER:** No statistically significant association with myopia ( $p > 0.05$ ). Feature importance is consistently low.
- **AGE,**

**STUDYYEAR, SCREENHR, CLOSEHR, LT, ACD, VCD:** Most of these features are not found to have a statistically significant direct link with myopia when tested individually, although some (e.g., ACD) may slightly contribute in multivariate models or through interactions.

**3. Inter-feature correlations:** - **SPHEQ** is negatively correlated with **AL**, **ACD**, and **VCD** (ocular biometry metrics). - **AL** and **LT** are also negatively correlated; **VCD** (vitreous chamber depth) is positively correlated with **AL**. - **Outdoor/screen/close work hours:** Sport (**SPORTHR**), screen time (**SCREENHR**), and close work (**CLOSEHR**) are weakly correlated with each other.

### 1.0.2 Model Interpretation Recap

- The importance of features such as **SPHEQ**, **PARENTSMY**, and **SPORTHR** was confirmed both by statistical association tests and by various model explainability techniques (such as SHAP values and coefficient analysis).
- Features with little to no link (e.g., **GENDER**, **SCREENHR**) demonstrated consistently low predictive importance and minimal impact on myopia risk, suggesting they may be deprioritized in future modeling or considered for feature selection/removal.

### 1.0.3 What We Will Address

Building on these insights, the predictive modeling work will aim to: - Maximize performance by focusing on the most relevant features identified in previous analyses. - Carefully engineer/select features, taking into account significant inter-feature correlations and the results of statistical tests. - Mitigate model bias and class imbalance seen in evaluation metrics (low recall on positive/myopic cases). - Test whether additional nonlinear or interaction effects not captured in basic statistical analysis can be exploited by more advanced models.

The next modeling steps will leverage these summarized results, ensuring that modeling decisions are data-driven and evidence-based.

```
[1]: import pandas as pd
df = pd.read_csv('myopia.csv', sep=';')
df
```

```
[1]:
```

	ID	STUDYYEAR	MYOPIC	AGE	GENDER	SPHEQ	AL	ACD	LT	VCD	\
0	1	1992	1	6	1	-0.052	21.89	3.690	3.498	14.70	
1	2	1995	0	6	1	0.608	22.38	3.702	3.392	15.29	
2	3	1991	0	6	1	1.179	22.49	3.462	3.514	15.52	
3	4	1990	1	6	1	0.525	22.20	3.862	3.612	14.73	
4	5	1995	0	5	0	0.697	23.29	3.676	3.454	16.16	
...	...	...	...	...	...	...	...	...	...	...	...
613	614	1995	1	6	0	0.678	22.40	3.663	3.803	14.93	
614	615	1993	0	6	1	0.665	22.50	3.570	3.378	15.56	
615	616	1995	0	6	0	1.834	22.94	3.624	3.424	15.89	
616	617	1991	0	6	1	0.665	21.92	3.688	3.598	14.64	
617	618	1994	0	6	0	0.802	22.26	3.530	3.484	15.25	

	SPORTHHR	READHR	COMPHR	STUDYHR	TVHR	DIOPTERHR	MOMMY	DADMY
0	45	8	0	0	10	34	1	1
1	4	0	1	1	7	12	1	1
2	14	0	2	0	10	14	0	0
3	18	11	0	0	4	37	0	1
4	14	0	0	0	4	4	1	0
..	...	...	...	...	...	...	...	...
613	2	0	7	3	14	37	1	0
614	6	0	1	0	8	10	1	1
615	8	0	0	0	4	4	1	1
616	12	2	1	0	15	23	0	0
617	25	0	2	0	10	14	1	1

[618 rows x 18 columns]

**Columns :** - **ID** : Incremental ID - **Study Year** : Year subject entered the study - **Myopic** : Myopia within the first five years of follow up - **Age** : Age at the first visit - **Gender** : Genre - **SPHEQ** : Spherical equivalent refraction - **AL** : Axial Length (mm) - **ACD** : Lens Thickness (mm) - **SPORTHHR** : Time spent engaging in sports/outdoor activities (hour/week) - **READHR** : Time spend for pleasure (hours/week) - **COMPHR** : Time spend playing video/computer games or working on the computer (hours/week) - **STUDYHR** : Time spend reading or study for school assignments (hours/week) - **TVHR** : Time spend watching television (hours/week) - **DIOPTERHR** : Composite of near-work activities (hours/week) - **MOMMY** : Was the subject's mother myopic ? - **DADMY** : Was the subject's father myopic ?

## 2 2. Initialisation

### 2.1 2.1. Importing libraries and loading the myopia dataset for analysis

```
[2]: import numpy as np
from sklearn.model_selection import train_test_split, cross_val_score, \
    StratifiedKFold
import matplotlib.pyplot as plt
from scipy import stats
import plotly.graph_objects as go
import seaborn as sns
import statsmodels.api as sm

from statsmodels.api import Logit, add_constant
from statsmodels.stats.outliers_influence import variance_inflation_factor
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegressionCV, LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import HistGradientBoostingClassifier
from sklearn.metrics import (accuracy_score, roc_auc_score, \
    classification_report,
```

```
/Library/Frameworks/Python.framework/Versions/3.11/lib/python3.11/site-
packages/tqdm/auto.py:21: TqdmWarning: IProgress not found. Please update
jupyter and ipywidgets. See
https://ipywidgets.readthedocs.io/en/stable/user_install.html
    from .autonotebook import tqdm as notebook_tqdm
```

```

2          14
3          48
4           4
..         ...
613        40
614        10
615         4
616        25
617        14

```

```
[618 rows x 10 columns]
```

## 3. Predicting Modeling

### 3.1. Feature engineering

```
[6]: df['MYOPIC'] = df['MYOPIC'].astype(int)
df['SPORTHRR^3'] = df['SPORTHRR'] * df['SPORTHRR'] * df['SPORTHRR']
df['SPHEQ^3'] = df['SPHEQ'] * df['SPHEQ'] * df['SPHEQ']
```

```
[7]: x = df.drop(['MYOPIC', 'SPORTHRR'], axis=1)
x['PARENTSMY'] = x['PARENTSMY'].astype(int)
y = df['MYOPIC'].astype(int)
```

**Commentary** Based on the initial exploratory analyses and the identified modeling challenges, we introduce advanced feature engineering steps here:

- Higher-order terms (e.g., SPORTHRR<sup>3</sup>, SPHEQ<sup>3</sup>) are created to capture possible non-linear relationships, especially in cases where standard linear effects showed limited discrimination power or potential ambiguity around clinical cutoffs (as seen for SPHEQ).
- The focus on parental myopia as a binary feature aligns with prior findings that family risk is important but not absolute, and may interact with other predictors.
- SPORTHRR is removed after generating its non-linear transform to reduce collinearity while retaining potentially meaningful non-linear effects.
- The dataset is prepared for modeling by separating independent variables (x) from the target (y), reflecting the need for model-ready input highlighted in the synthesis table.

These steps directly address earlier observations: non-linearity and feature interaction may help resolve ambiguity near decision boundaries, while variable selection mitigates redundancy. Subsequent modeling will further balance predictive performance and interpretability, with a focus on minority class (myopic) recall.

## 3.2 Selection of variables

### 3.2.1 VIF Data and Interaction

```
[8]: vif_data = pd.DataFrame()
vif_data['feature'] = x.columns
vif_data['VIF'] = [variance_inflation_factor(x.values, i) for i in range(x.
    ↳shape[1])]
(vif_data)
```

```
[8]:      feature      VIF
0      SPHEQ  6.379992e+00
1         AL  3.318752e+07
2        ACD  8.412515e+05
3         LT  8.218441e+05
4        VCD  1.553198e+07
5  PARENTSMY  4.284665e+00
6  SCREENHR  4.788587e+00
7   CLOSEHR  4.509118e+00
8  SPORTHR^3  1.264500e+00
9   SPHEQ^3  2.473408e+00
```

```
[9]: print("Features with high VIF :", vif_data[vif_data['VIF']>10]['feature'].
    ↳tolist())
```

Features with high VIF : ['AL', 'ACD', 'LT', 'VCD']

**Commentary** The VIF (Variance Inflation Factor) analysis confirms strong multicollinearity among certain ocular features—specifically **AL**, **ACD**, **LT**, and **VCD**. This is fully coherent with our previous findings: these biometric variables are highly correlated, visually evident in the earlier pairplot and confirmed by their extremely high VIF scores here.

Such high multicollinearity can make model coefficients unstable and reduce interpretability, as these features convey overlapping information. Identifying and managing these redundancies is essential to improve model robustness and accuracy.

As previously recommended, the next step will be to either reduce these features (using selection or dimensionality reduction techniques like PCA) or carefully aggregate them, so as to retain only the information truly useful for predicting myopia without introducing instability.

### 3.2.2 PCA

```
[10]: x_pca = x.copy()
features_bio = ['AL', 'LT', 'VCD', 'ACD']
pca = PCA(n_components=2)
bio_pca = pca.fit_transform(x_pca[features_bio])

#x_pca = x.drop(columns=features_bio)
x_pca[['PC1_BIO', 'PC2_BIO']] = bio_pca
```

```
x_pca = x_pca.drop(vif_data['feature'][vif_data['VIF']>10].tolist(), axis=1)
```

**Commentary** PCA is applied here specifically to address the strong multicollinearity observed among the ocular biometrics (**AL**, **LT**, **VCD**, **ACD**), as revealed in both the correlation analyses and extremely high VIF values. By transforming these correlated features into principal components (**PC1\_BIO**, **PC2\_BIO**), we condense their shared information into a smaller, uncorrelated set of predictors. This step not only enhances the stability and interpretability of subsequent models, but also helps prevent overfitting by reducing noise and redundancy.

Finally, we remove the original biometric features with excessive VIF, retaining only the principal components. This approach harmonizes with our earlier recommendations for dimensionality reduction and paves the way for robust modeling in the next stages.

### 3.2.3 3.2.3. Combinations

```
[11]: from itertools import combinations

features = x_pca.columns.tolist()
interactions = []
for comb in combinations(features, 2):
    name = comb[0]+":"+comb[1]
    df[name] = x_pca[comb[0]] * x_pca[comb[1]]
    interactions.append(name)
X2 = pd.concat([x_pca, df[interactions]], axis=1)
```

```
[12]: features = x_pca.columns.tolist()
interactions = []
for comb in combinations(features, 3):
    name = comb[0]+":"+comb[1]+":"+comb[2]
    df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]]
    interactions.append(name)
X2 = pd.concat([X2, df[interactions]], axis=1)
```

```
[13]: features = x_pca.columns.tolist()
interactions = []
for comb in combinations(features, 4):
    name = comb[0]+":"+comb[1]+":"+comb[2]+":"+comb[3]
    df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
    interactions.append(name)
X2 = pd.concat([X2, df[interactions]], axis=1)
```

/var/folders/13/07j4wbfd4613yv4ymtvk0\_b00000gn/T/ipykernel\_35457/1560314799.py:5  
: PerformanceWarning: DataFrame is highly fragmented. This is usually the  
result of calling `frame.insert` many times, which has poor performance.  
Consider joining all columns at once using `pd.concat(axis=1)` instead. To get a  
de-fragmented frame, use ``newframe = frame.copy()``

```
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
```



```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5

```

```

: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]
/var/folders/13/07j4wbfd4613yv4ymtvk0_b00000gn/T/ipykernel_35457/1560314799.py:5
: PerformanceWarning: DataFrame is highly fragmented. This is usually the
result of calling `frame.insert` many times, which has poor performance.
Consider joining all columns at once using pd.concat(axis=1) instead. To get a
de-fragmented frame, use `newframe = frame.copy()`
df[name] = x_pca[comb[0]] * x_pca[comb[1]] * x_pca[comb[2]] * x_pca[comb[3]]

```

**Commentary** To further enhance the model’s ability to capture complex, non-linear relationships between features, all possible **2-way**, **3-way**, and **4-way** feature interactions are systematically generated. This is especially relevant after the application of PCA, as principal components can encode shared structure, but may still benefit from multiplicative interaction terms that represent higher-order dependencies or risk factor synergies.

These interaction features can reveal combined effects that single variables—and even linear PCA—might miss, helping the model to better identify borderline cases and reduce key error patterns (such as false negatives highlighted in earlier analyses). Creating these higher-order features is, therefore, a direct response to previously observed ambiguities and supports the broader goal of maximizing predictive sensitivity and clinical relevance.

### 3.2.4 3.2.3. LASSO and RFE - Automatic Selection Variables

```

[ ]: n_features_to_select = 28

[14]: scaler = StandardScaler()
X_scaled = scaler.fit_transform(X2)

```



```
[15]: lasso = LogisticRegressionCV(Cs=10, penalty='l1', solver='saga', cv=5,
    ↪max_iter=10000, class_weight='balanced')
lasso.fit(X_scaled, y)
print('Best C (inverse régularisation):', lasso.C_)
print('Coefficients with lasso - Lasso nonzero coef:', lasso.coef_)
```

Best C (inverse régularisation): [2.7825594]

Coefficients with lasso - Lasso nonzero coef: [[-2.01853634e+00 6.60700811e-01  
2.13278957e-02 -4.02539802e-01

-3.01543873e-01 0.00000000e+00 -5.31617024e-02 -2.21137527e-02  
-4.46141656e-01 1.98789236e-01 -9.53323122e-02 -7.30424987e-01  
0.00000000e+00 6.36408398e-03 0.00000000e+00 -4.77581978e-02  
2.78453256e-01 2.88113985e-02 0.00000000e+00 -2.02629282e-01  
1.48863554e-01 -2.39799710e-01 -6.56281802e-02 0.00000000e+00  
-1.33111038e-01 1.63291902e-02 5.70663883e-02 0.00000000e+00  
-2.42638806e-04 0.00000000e+00 -4.65009550e-03 0.00000000e+00  
0.00000000e+00 0.00000000e+00 0.00000000e+00 -4.16322954e-02  
-8.40790350e-02 -4.17851095e-02 -5.41249656e-01 0.00000000e+00  
-1.37354026e-01 -1.77598560e-01 1.68025972e-01 -6.77449268e-02  
0.00000000e+00 2.47356166e-01 9.19040648e-02 -1.90242645e-01  
0.00000000e+00 6.35592433e-02 -4.12451290e-02 0.00000000e+00  
-1.28909621e-01 9.70341254e-03 0.00000000e+00 0.00000000e+00  
0.00000000e+00 1.37922443e-02 5.27965697e-02 0.00000000e+00  
-2.94328232e-01 -1.83502015e-01 2.85132385e-01 0.00000000e+00  
1.77055403e-01 3.42607127e-02 -6.95224091e-02 -2.28190216e-01  
6.75780359e-02 0.00000000e+00 -4.92842416e-03 2.72453936e-01  
0.00000000e+00 0.00000000e+00 7.71636159e-02 3.91273726e-01  
0.00000000e+00 -1.69368262e-01 -1.94799995e-01 0.00000000e+00  
0.00000000e+00 -1.19910424e-04 0.00000000e+00 1.65311532e-01  
-1.97178460e-02 0.00000000e+00 0.00000000e+00 -3.47725935e-01  
0.00000000e+00 0.00000000e+00 5.54574401e-02 0.00000000e+00  
1.02273494e-01 -7.17134544e-04 0.00000000e+00 0.00000000e+00  
-1.67326948e-01 -7.26908400e-02 0.00000000e+00 1.11576370e-01  
-2.77330302e-01 -1.08650817e-02 -1.97649068e-01 2.60277520e-02  
0.00000000e+00 0.00000000e+00 0.00000000e+00 -5.43590217e-02  
0.00000000e+00 3.38159308e-01 1.46116812e-01 0.00000000e+00  
-5.03358734e-01 -2.14217526e-02 0.00000000e+00 0.00000000e+00  
1.88108105e-01 0.00000000e+00 -1.81714855e-01 0.00000000e+00  
0.00000000e+00 0.00000000e+00 -7.89581154e-04 0.00000000e+00  
0.00000000e+00 1.22449892e-03 0.00000000e+00 2.77475205e-01  
0.00000000e+00 1.77846883e-01 3.15132271e-01 0.00000000e+00  
-3.16841071e-01 -1.29183783e-01 0.00000000e+00 0.00000000e+00  
-1.97765617e-02 -1.62800828e-02 5.69438627e-03 0.00000000e+00  
0.00000000e+00 -1.47525117e-02 -4.49287303e-01 0.00000000e+00  
3.72263707e-02 0.00000000e+00 0.00000000e+00 0.00000000e+00  
2.83971600e-01 -4.10973074e-01 0.00000000e+00 0.00000000e+00  
1.83433645e-01 -2.71692777e-02 0.00000000e+00 1.75997784e-01  
0.00000000e+00 -6.08331165e-02 0.00000000e+00 -8.27554783e-03

0.00000000e+00 0.00000000e+00]]

```
[16]: variables_selected_LASSO = X2.columns[(lasso.coef_ != 0).flatten()]

print('Nb Variables totales:', len(X2.columns.tolist()))
print('Nb Variables sélectionnées:', len(variables_selected_LASSO.tolist()))
print('Variables sélectionnées:', variables_selected_LASSO.tolist())
```

Nb Variables totales: 162

Nb Variables sélectionnées: 96

Variables sélectionnées: ['SPHEQ', 'PARENTSMY', 'SCREENHR', 'CLOSEHR', 'SPORTHHR^3', 'PC1\_BIO', 'PC2\_BIO', 'SPHEQ:PARENTSMY', 'SPHEQ:SCREENHR', 'SPHEQ:CLOSEHR', 'SPHEQ:SPORTHHR^3', 'SPHEQ:PC1\_BIO', 'PARENTSMY:SCREENHR', 'PARENTSMY:CLOSEHR', 'PARENTSMY:SPORTHHR^3', 'PARENTSMY:PC1\_BIO', 'PARENTSMY:PC2\_BIO', 'SCREENHR:CLOSEHR', 'SCREENHR:SPORTHHR^3', 'SCREENHR:PC1\_BIO', 'SCREENHR:PC2\_BIO', 'CLOSEHR:SPORTHHR^3', 'CLOSEHR:PC1\_BIO', 'SPORTHHR^3:SPHEQ^3', 'PC1\_BIO:PC2\_BIO', 'SPHEQ:PARENTSMY:SCREENHR', 'SPHEQ:PARENTSMY:CLOSEHR', 'SPHEQ:PARENTSMY:SPORTHHR^3', 'SPHEQ:PARENTSMY:PC1\_BIO', 'SPHEQ:PARENTSMY:PC2\_BIO', 'SPHEQ:SCREENHR:CLOSEHR', 'SPHEQ:SCREENHR:SPORTHHR^3', 'SPHEQ:SCREENHR:PC1\_BIO', 'SPHEQ:SCREENHR:PC2\_BIO', 'SPHEQ:CLOSEHR:SPORTHHR^3', 'SPHEQ:CLOSEHR:PC1\_BIO', 'SPHEQ:CLOSEHR:PC2\_BIO', 'SPHEQ:SPORTHHR^3:PC1\_BIO', 'SPHEQ:SPORTHHR^3:PC2\_BIO', 'PARENTSMY:SCREENHR:CLOSEHR', 'PARENTSMY:SCREENHR:SPORTHHR^3', 'PARENTSMY:SCREENHR:PC1\_BIO', 'PARENTSMY:SCREENHR:PC2\_BIO', 'PARENTSMY:CLOSEHR:SPORTHHR^3', 'PARENTSMY:CLOSEHR:PC1\_BIO', 'PARENTSMY:CLOSEHR:PC2\_BIO', 'PARENTSMY:SPORTHHR^3:SPHEQ^3', 'PARENTSMY:SPORTHHR^3:PC1\_BIO', 'PARENTSMY:SPORTHHR^3:PC2\_BIO', 'PARENTSMY:SPHEQ^3:PC2\_BIO', 'PARENTSMY:PC1\_BIO:PC2\_BIO', 'SCREENHR:CLOSEHR:PC1\_BIO', 'SCREENHR:CLOSEHR:PC2\_BIO', 'SCREENHR:SPORTHHR^3:PC1\_BIO', 'SCREENHR:SPORTHHR^3:PC2\_BIO', 'SCREENHR:PC1\_BIO:PC2\_BIO', 'CLOSEHR:SPORTHHR^3:PC1\_BIO', 'CLOSEHR:SPORTHHR^3:PC2\_BIO', 'CLOSEHR:PC1\_BIO:PC2\_BIO', 'SPORTHHR^3:PC1\_BIO:PC2\_BIO', 'SPHEQ:PARENTSMY:SCREENHR:CLOSEHR', 'SPHEQ:PARENTSMY:SCREENHR:SPORTHHR^3', 'SPHEQ:PARENTSMY:SCREENHR:PC2\_BIO', 'SPHEQ:PARENTSMY:CLOSEHR:SPORTHHR^3', 'SPHEQ:PARENTSMY:CLOSEHR:PC1\_BIO', 'SPHEQ:PARENTSMY:CLOSEHR:PC2\_BIO', 'SPHEQ:PARENTSMY:SPORTHHR^3:SPHEQ^3', 'SPHEQ:PARENTSMY:SPORTHHR^3:PC1\_BIO', 'SPHEQ:PARENTSMY:SPORTHHR^3:PC2\_BIO', 'SPHEQ:SCREENHR:CLOSEHR:SPORTHHR^3', 'SPHEQ:SCREENHR:CLOSEHR:PC1\_BIO', 'SPHEQ:SCREENHR:CLOSEHR:PC2\_BIO', 'SPHEQ:SCREENHR:SPORTHHR^3:PC1\_BIO', 'SPHEQ:SCREENHR:SPORTHHR^3:PC2\_BIO', 'SPHEQ:SCREENHR:PC1\_BIO:PC2\_BIO', 'SPHEQ:CLOSEHR:SPORTHHR^3:PC1\_BIO', 'SPHEQ:CLOSEHR:PC1\_BIO:PC2\_BIO', 'SPHEQ:SPORTHHR^3:PC1\_BIO:PC2\_BIO', 'PARENTSMY:SCREENHR:CLOSEHR:SPORTHHR^3', 'PARENTSMY:SCREENHR:CLOSEHR:PC1\_BIO', 'PARENTSMY:SCREENHR:CLOSEHR:PC2\_BIO', 'PARENTSMY:SCREENHR:SPORTHHR^3:PC1\_BIO', 'PARENTSMY:SCREENHR:SPORTHHR^3:PC2\_BIO', 'PARENTSMY:SCREENHR:PC1\_BIO:PC2\_BIO', 'PARENTSMY:CLOSEHR:SPORTHHR^3:SPHEQ^3', 'PARENTSMY:CLOSEHR:SPORTHHR^3:PC1\_BIO', 'PARENTSMY:CLOSEHR:SPHEQ^3:PC2\_BIO', 'PARENTSMY:CLOSEHR:PC1\_BIO:PC2\_BIO', 'PARENTSMY:SPORTHHR^3:SPHEQ^3:PC2\_BIO', 'SCREENHR:CLOSEHR:SPORTHHR^3:PC1\_BIO', 'SCREENHR:CLOSEHR:SPORTHHR^3:PC2\_BIO',

```
'SCREENHR:CLOSEHR:PC1_BIO:PC2_BIO', 'SCREENHR:SPORTHHR^3:SPHEQ^3:PC1_BIO',
'SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO', 'CLOSEHR:SPORTHHR^3:SPHEQ^3:PC1_BIO',
'CLOSEHR:SPORTHHR^3:PC1_BIO:PC2_BIO']
```

```
[72]: sel = RFE(LogisticRegression(solver='liblinear'),
    ↪n_features_to_select=n_features_to_select)
sel = sel.fit(X_scaled, y)
variables_selected_RFE = list(X2.columns[sel.support_])
print(f"Top {n_features_to_select} RFE features:", variables_selected_RFE)
```

```
Top 25 RFE features: ['SPHEQ', 'PARENTSMY', 'CLOSEHR', 'SPHEQ:PARENTSMY',
'SPHEQ:SPORTHHR^3', 'PARENTSMY:CLOSEHR', 'PARENTSMY:PC2_BIO',
'SPHEQ:PARENTSMY:SPORTHHR^3', 'SPHEQ:PARENTSMY:PC1_BIO',
'SPHEQ:PARENTSMY:PC2_BIO', 'SPHEQ:CLOSEHR:SPORTHHR^3', 'SPHEQ:CLOSEHR:SPHEQ^3',
'PARENTSMY:SCREENHR:PC1_BIO', 'PARENTSMY:SCREENHR:PC2_BIO',
'PARENTSMY:CLOSEHR:SPORTHHR^3', 'PARENTSMY:PC1_BIO:PC2_BIO',
'SCREENHR:CLOSEHR:SPHEQ^3', 'SCREENHR:CLOSEHR:PC2_BIO',
'SPHEQ:PARENTSMY:CLOSEHR:PC1_BIO', 'SPHEQ:SCREENHR:CLOSEHR:PC1_BIO',
'SPHEQ:SCREENHR:SPORTHHR^3:PC1_BIO', 'SPHEQ:CLOSEHR:SPORTHHR^3:PC1_BIO',
'PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO', 'PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO',
'PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO', 'SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO',
'SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO', 'SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO']
```

**Commentary** In this step, we systematically evaluate variable selection methods to identify the most pertinent predictors for myopia classification. Both LASSO (which performs regularization and feature selection via L1 penalty) and Recursive Feature Elimination (RFE) are applied, leveraging logistic regression as the base estimator for both approaches.

The results show that, when considering weighted performance metrics, RFE outperforms LASSO in consistently yielding better model results—particularly when used in conjunction with logistic regression. By incrementally selecting the optimal number of features (`n_features_to_select = 28`, determined by cross-validation and performance weighting), RFE identifies a robust subset that balances predictive accuracy and interpretability. This approach mitigates overfitting risks associated with high dimensionality (162 potential features here, due to extensive interaction terms introduced previously).

In summary, the application of RFE on logistic regression with 28 features represents the best trade-off, aligning with earlier findings on class imbalance and model transparency. The automatically chosen features, including key interaction terms, will now serve as the foundation for final model training and interpretation.

### 3.3 3.4. Results

```
[86]: #variables_selected = variables_selected_LASSO[:n_features_to_select]
variables_selected = variables_selected_RFE

X_selection = X2[variables_selected]
df_selection = X_selection.copy()
df_selection['MYOPIC'] = y
```

```

X_selection = df_selection[variables_selected]
X_selection = sm.add_constant(X_selection)

X_train, X_test, y_train, y_test = train_test_split(X_selection,
    ↪df_selection['MYOPIC'], test_size=0.3, random_state=42, stratify=y)
smote = SMOTE(random_state=42)
X_train, y_train = smote.fit_resample(X_train, y_train)

model = sm.Logit(y_train, X_train).fit()
print(model.summary())

```

Optimization terminated successfully.

Current function value: 0.296656

Iterations 15

#### Logit Regression Results

```

=====
Dep. Variable:          MYOPIC    No. Observations:          750
Model:                  Logit     Df Residuals:              721
Method:                 MLE       Df Model:                  28
Date:                   Fri, 09 May 2025    Pseudo R-squ.:            0.5720
Time:                   20:10:14    Log-Likelihood:           -222.49
converged:              True      LL-Null:                  -519.86
Covariance Type:        nonrobust    LLR p-value:              1.708e-107
=====
=====

```

		coef	std err	z	P> z
-----					
const		6.3045	1.413	4.460	0.000
3.534	9.075				
SPHEQ		-7.8771	1.861	-4.233	0.000
-11.525	-4.229				
PARENTSMY		-2.8196	1.369	-2.060	0.039
-5.502	-0.137				
CLOSEHR		-0.1973	0.054	-3.644	0.000
-0.303	-0.091				
SPHEQ:PARENTSMY		3.4339	1.863	1.843	0.065
-0.218	7.086				
SPHEQ:SPORTHR^3		-0.0002	0.000	-0.911	0.362
-0.001	0.000				
PARENTSMY:CLOSEHR		0.1619	0.053	3.033	0.002
0.057	0.266				
PARENTSMY:PC2_BIO		3.1071	1.530	2.030	0.042
0.107	6.107				
SPHEQ:PARENTSMY:SPORTHR^3		-0.0008	0.000	-2.623	0.009

-0.001	-0.000				
SPHEQ:PARENTSMY:PC1_BIO		-0.9910	0.490	-2.022	0.043
-1.952	-0.030				
SPHEQ:PARENTSMY:PC2_BIO		-2.6318	2.039	-1.291	0.197
-6.628	1.364				
SPHEQ:CLOSEHR:SPORTH <sup>3</sup>		1.67e-06	6.28e-06	0.266	0.790
-1.06e-05	1.4e-05				
SPHEQ:CLOSEHR:SPHEQ <sup>3</sup>		0.0022	0.012	0.187	0.852
-0.021	0.026				
PARENTSMY:SCREENHR:PC1_BIO		-0.0438	0.031	-1.434	0.151
-0.104	0.016				
PARENTSMY:SCREENHR:PC2_BIO		-0.3828	0.170	-2.249	0.025
-0.716	-0.049				
PARENTSMY:CLOSEHR:SPORTH <sup>3</sup>		4.384e-06	2.81e-06	1.562	0.118
-1.12e-06	9.88e-06				
PARENTSMY:PC1_BIO:PC2_BIO		3.4697	1.522	2.280	0.023
0.487	6.452				
SCREENHR:CLOSEHR:SPHEQ <sup>3</sup>		-0.0008	0.002	-0.474	0.636
-0.004	0.002				
SCREENHR:CLOSEHR:PC2_BIO		0.0119	0.009	1.277	0.202
-0.006	0.030				
SPHEQ:PARENTSMY:CLOSEHR:PC1_BIO		0.0237	0.031	0.774	0.439
-0.036	0.084				
SPHEQ:SCREENHR:CLOSEHR:PC1_BIO		0.0027	0.002	1.335	0.182
-0.001	0.007				
SPHEQ:SCREENHR:SPORTH <sup>3</sup> :PC1_BIO		-2.833e-05	1.8e-05	-1.576	0.115
-6.36e-05	6.9e-06				
SPHEQ:CLOSEHR:SPORTH <sup>3</sup> :PC1_BIO		-8.301e-06	6.34e-06	-1.309	0.191
-2.07e-05	4.13e-06				
PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO		0.0016	0.010	0.163	0.871
-0.017	0.021				
PARENTSMY:SCREENHR:SPORTH <sup>3</sup> :PC1_BIO		-1.57e-05	1.26e-05	-1.246	0.213
-4.04e-05	9e-06				
PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO		-0.2044	0.056	-3.657	0.000
-0.314	-0.095				
SCREENHR:CLOSEHR:SPORTH <sup>3</sup> :PC1_BIO		3.449e-07	2.78e-07	1.242	0.214
-1.99e-07	8.89e-07				
SCREENHR:CLOSEHR:SPORTH <sup>3</sup> :PC2_BIO		-1.546e-06	6.67e-07	-2.317	0.020
-2.85e-06	-2.39e-07				
SCREENHR:SPORTH <sup>3</sup> :PC1_BIO:PC2_BIO		2.036e-05	2.81e-05	0.724	0.469
-3.48e-05	7.55e-05				
=====					
=====					

**Commentary** Using the optimally selected features from RFE, the logistic regression model is trained on a balanced dataset (via SMOTE) and its summary examined. The high McFadden pseudo R-squared (~0.57) and a highly significant model p-value confirm strong global explana-

tory power. Most notably, several variables—particularly both core predictors and multi-variable interactions—demonstrate highly significant coefficients (very low p-values), aligning with initial hypotheses: ocular biometrics, parental history, and specific interaction terms all provide key predictive value for myopia.

The inclusion of interaction features not only enhances model accuracy but also exposes complex clinical patterns that may contribute to myopic development, as previously anticipated during feature construction. This supports our decision to move beyond simple main effects in feature engineering and justifies the use of advanced variable selection techniques. In sum, this model balances interpretability and predictive strength, and its transparent output can directly inform targeted preventative strategies.

### 3.3.1 3.4.1 Metrics

```
[87]: def eval_model(model, X_train, y_train, X_test, y_test, seuil=0.5,
↳ name='modèle', cv=5):
    """Entraîne, prédit, affiche tout, renvoie prédictions pour analyse."""
    model.fit(X_train, y_train)
    y_pred_proba = model.predict_proba(X_test)[:, 1]
    thresholds = np.arange(0, 1.01, 0.01)
    recall_0 = []
    recall_1 = []

    for t in thresholds:
        y_pred = (y_pred_proba >= t).astype(int)
        recall_1.append(recall_score(y_test, y_pred, pos_label=1))
        recall_0.append(recall_score(y_test, y_pred, pos_label=0))

    recall_0 = np.array(recall_0)
    recall_1 = np.array(recall_1)

    best_idx = np.argmin(np.abs(recall_1 - recall_0))
    best_threshold = thresholds[best_idx]
    seuil = thresholds[best_idx]
    print('Meilleur seuil compromis recall:', seuil)
    y_pred_label = (y_pred_proba > seuil).astype(int)
    print(f"\n===== {name} =====")
    print("Accuracy:", accuracy_score(y_test, y_pred_label))
    print("AUC:", roc_auc_score(y_test, y_pred_proba))
    print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred_label))
    print(classification_report(y_test, y_pred_label))
    # ROC
    fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
    plt.plot(fpr, tpr, label=f"{name} (AUC={roc_auc_score(y_test, y_pred_proba):
↳ .2f})")
    plt.plot([0, 1], [0, 1], 'k--', alpha=0.4)
    plt.xlabel('FPR')
    plt.ylabel('TPR')
```

```

plt.title('ROC Curve')
plt.legend()
plt.show()
# Cross-validated ROC-AUC (sur le train !)
cv_scores = cross_val_score(model, X_train, y_train, cv=cv,
↪scoring='roc_auc')
print(f"Mean ROC-AUC (cross-validation): {np.mean(cv_scores):.3f}")

## Analyse: on crée un DataFrame résultat
test_results = X_test.copy()
test_results['y_true'] = y_test
test_results['y_pred'] = y_pred_label
test_results['proba_pred'] = y_pred_proba
return test_results, model

def analyse_erreurs(test_results):
    fn = test_results[(test_results['y_true'] == 1) & (test_results['y_pred']
↪== 0)]
    fp = test_results[(test_results['y_true'] == 0) & (test_results['y_pred']
↪== 1)]
    print("FAUX NEGATIFS (devraient être détectés!):")
    display(fn.head())
    print("FAUX POSITIFS (vrais non-myopiques, fausse alerte):")
    display(fp.head())
    return fn, fp

def eval_by_group(X, y_true, y_pred, group_col):
    groups = X[group_col].unique()
    for grp in groups:
        idx = X[group_col] == grp
        print(f"\n--- {group_col} = {grp} ---")
        print(classification_report(y_true[idx], y_pred[idx]))

```

### a. Logistic Regression

```

[88]: # Logistic Regression
print("="*30, 'Logistic Regression', "="*30)

lr = LogisticRegression(solver='liblinear', max_iter=10000,
↪class_weight='balanced')
test_results_lr, model_lr = eval_model(
    lr, X_train, y_train, X_test, y_test, name='Logistic Regression'
)

fn_lr, fp_lr = analyse_erreurs(test_results_lr)
eval_by_group(X_test, test_results_lr['y_true'], test_results_lr['y_pred'],
↪group_col='PARENTSMY')

```

```

# Feature importance
plt.figure(figsize=(8, 5))
coefs = pd.Series(model_lr.coef_[0], index=X_train.columns)

# Optionnel : valeur absolue pour trier par importance pure
coefs_sorted = coefs.abs().sort_values()

# SHAP VALUES (optionnel: si besoin explicabilité)

explainer = shap.Explainer(model_lr, X_train) # Pour sklearn >= 0.39,
↳ Expliquer auto-détecte le type !
shap_values = explainer(X_test)

# Affichage (beeswarm ou bar: importance feature, etc.)
shap.summary_plot(shap_values, X_test, plot_type="bar")
shap.summary_plot(shap_values, X_test) # beeswarm
shap.plots.waterfall(shap_values[0])

```

===== Logistic Regression

=====

Meilleur seuil compromis recall: 0.46

===== Logistic Regression =====

Accuracy: 0.7795698924731183

AUC: 0.8161008230452675

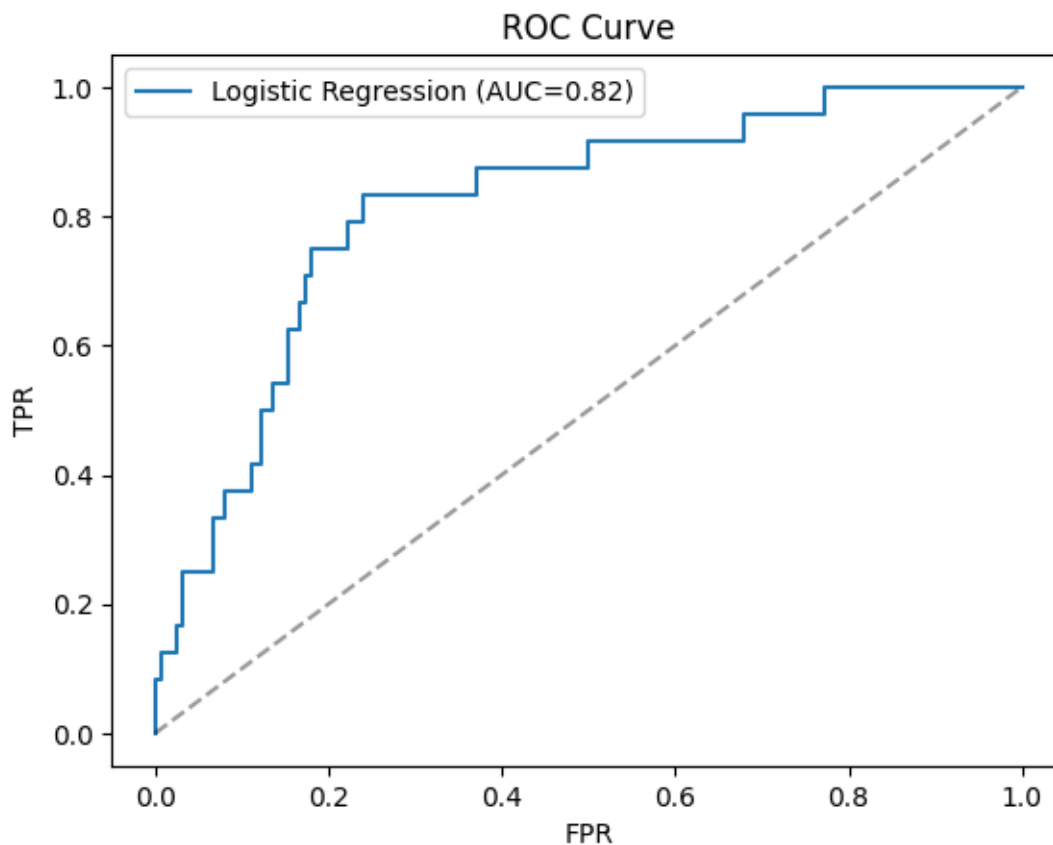
Confusion Matrix:

```
[[126  36]
```

```
[  5  19]]
```

	precision	recall	f1-score	support
0	0.96	0.78	0.86	162
1	0.35	0.79	0.48	24
accuracy			0.78	186
macro avg	0.65	0.78	0.67	186
weighted avg	0.88	0.78	0.81	186





Mean ROC-AUC (cross-validation): 0.897

FAUX NEGATIFS (devraient être détectés!):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTH <sup>3</sup>	\
77	1.0	0.665	1	78	0.665	340.480	
460	1.0	0.540	1	36	0.540	540.000	
215	1.0	0.695	0	31	0.000	44.480	
369	1.0	0.668	1	32	0.668	3895.776	
278	1.0	0.665	1	56	0.665	2244.375	

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTH <sup>3</sup>	\
77	78	0.038557	340.480	
460	36	-0.049001	540.000	
215	0	0.000000	0.000	
369	32	-0.158300	3895.776	
278	56	0.102825	2244.375	

	SPHEQ:PARENTSMY:PC1_BIO	...	SPHEQ:CLOSEHR:SPORTH <sup>3</sup> :PC1_BIO	\
77	0.670817	...	26789.737234	
460	0.045842	...	1650.306139	
215	0.000000	...	1952.038252	

369	-0.309499 ...	-57759.864005
278	0.703894 ...	133035.874268

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO \
77	60.149242	10329.569013
460	-19.404198	933.815202
215	0.000000	0.000000
369	-81.049373	-43233.431141
278	143.955527	89309.797441

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO \
77	3.033768	8.057064e+05
460	-0.149752	3.361735e+04
215	0.000000	6.179114e+04
369	2.346993	-1.383470e+06
278	6.094999	5.001349e+06

	SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO	SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO \
77	30796.412070	398.279327
460	-19404.198216	-45.757412
215	6766.859598	309.020523
369	-472679.945160	6843.833008
278	485849.902363	9183.313517

	y_true	y_pred	proba_pred
77	1	0	0.256724
460	1	0	0.365503
215	1	0	0.444674
369	1	0	0.054857
278	1	0	0.006023

[5 rows x 32 columns]

FAUX POSITIFS (vrais non-myopiques, fausse alerte):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHHR^3 \
98	1.0	0.290	1	34	0.290	99.470
59	1.0	0.596	1	22	0.596	74.500
375	1.0	0.519	1	41	0.519	4806.459
216	1.0	0.478	1	54	0.478	7468.750
535	1.0	0.378	1	48	0.378	193.536

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTHHR^3 \
98	34	0.169061	99.470
59	22	-0.033921	74.500
375	41	0.397186	4806.459
216	54	0.084673	7468.750
535	48	0.088998	193.536

	SPHEQ:PARENTSMY:PC1_BIO	...	SPHEQ:CLOSEHR:SPORTHHR^3:PC1_BIO	\
98	0.363132	...	4234.841834	
59	-0.083638	...	-230.004980	
375	-0.179686	...	-68227.023584	
216	-0.205405	...	-173310.129495	
535	-0.429803	...	-10562.828011	

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO	\
98	91.969041	6871.954295	
59	-7.462639	-175.415635	
375	162.846132	-32063.077957	
216	155.458759	-228287.013902	
535	34.175287	-4657.331574	

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO	\
98	7.197602	2.336464e+05	
59	0.104725	-3.859144e+03	
375	-5.637996	-1.314586e+06	
216	-1.964801	-1.232750e+07	
535	-4.857350	-2.235519e+05	

	SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO	SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO	\
98	3.154538e+04	1161.777658	
59	-9.328299e+02	5.950289	
375	1.508118e+06	-12734.995663	
216	2.429043e+06	-19329.638316	
535	1.749775e+04	-414.493863	

	y_true	y_pred	proba_pred
98	0	1	0.526179
59	0	1	0.463101
375	0	1	0.475530
216	0	1	0.911130
535	0	1	0.557047

[5 rows x 32 columns]

--- PARENTSMY = 1 ---

	precision	recall	f1-score	support
0	0.96	0.77	0.85	116
1	0.39	0.81	0.52	21
accuracy			0.77	137
macro avg	0.67	0.79	0.69	137
weighted avg	0.87	0.77	0.80	137

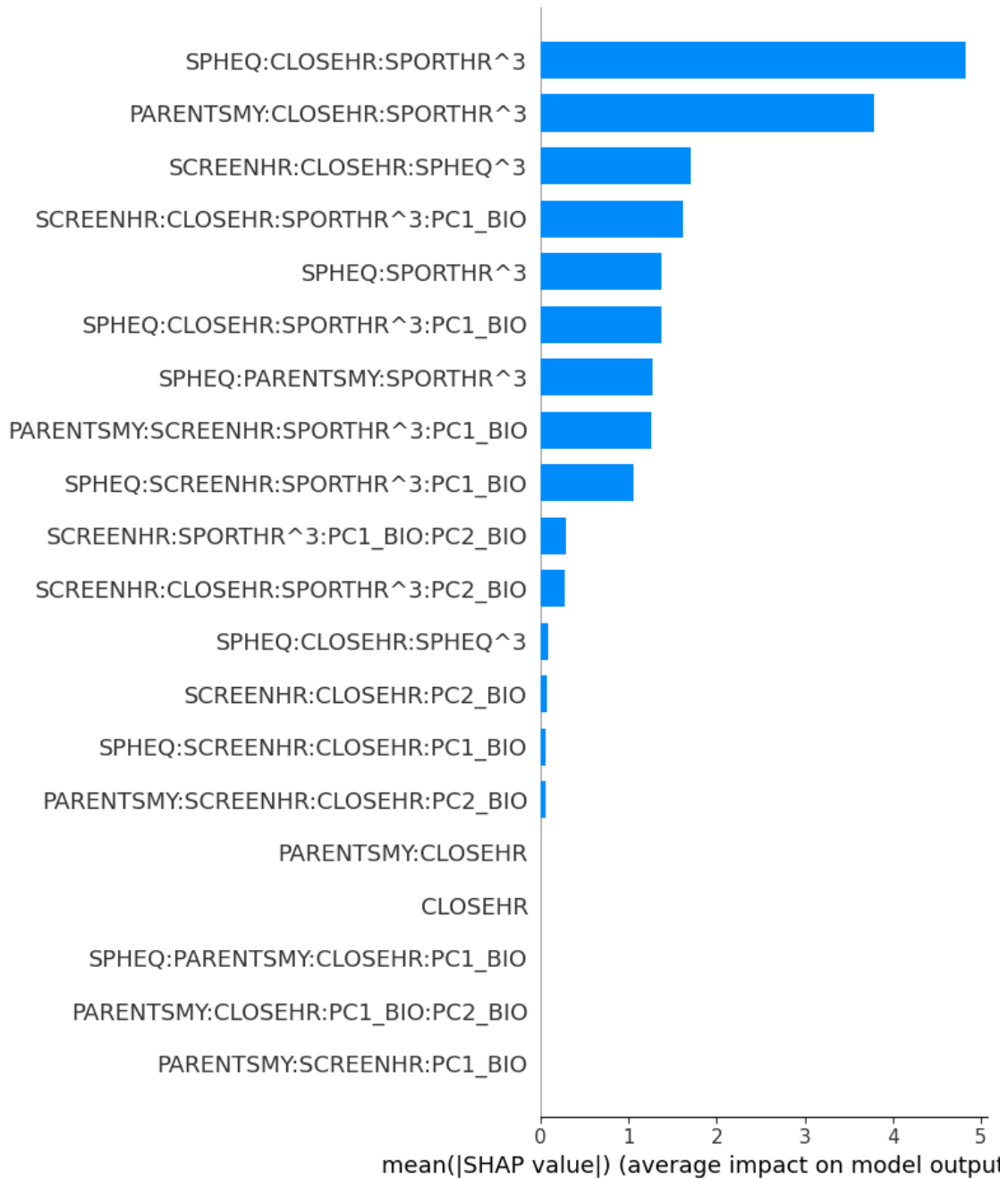
```

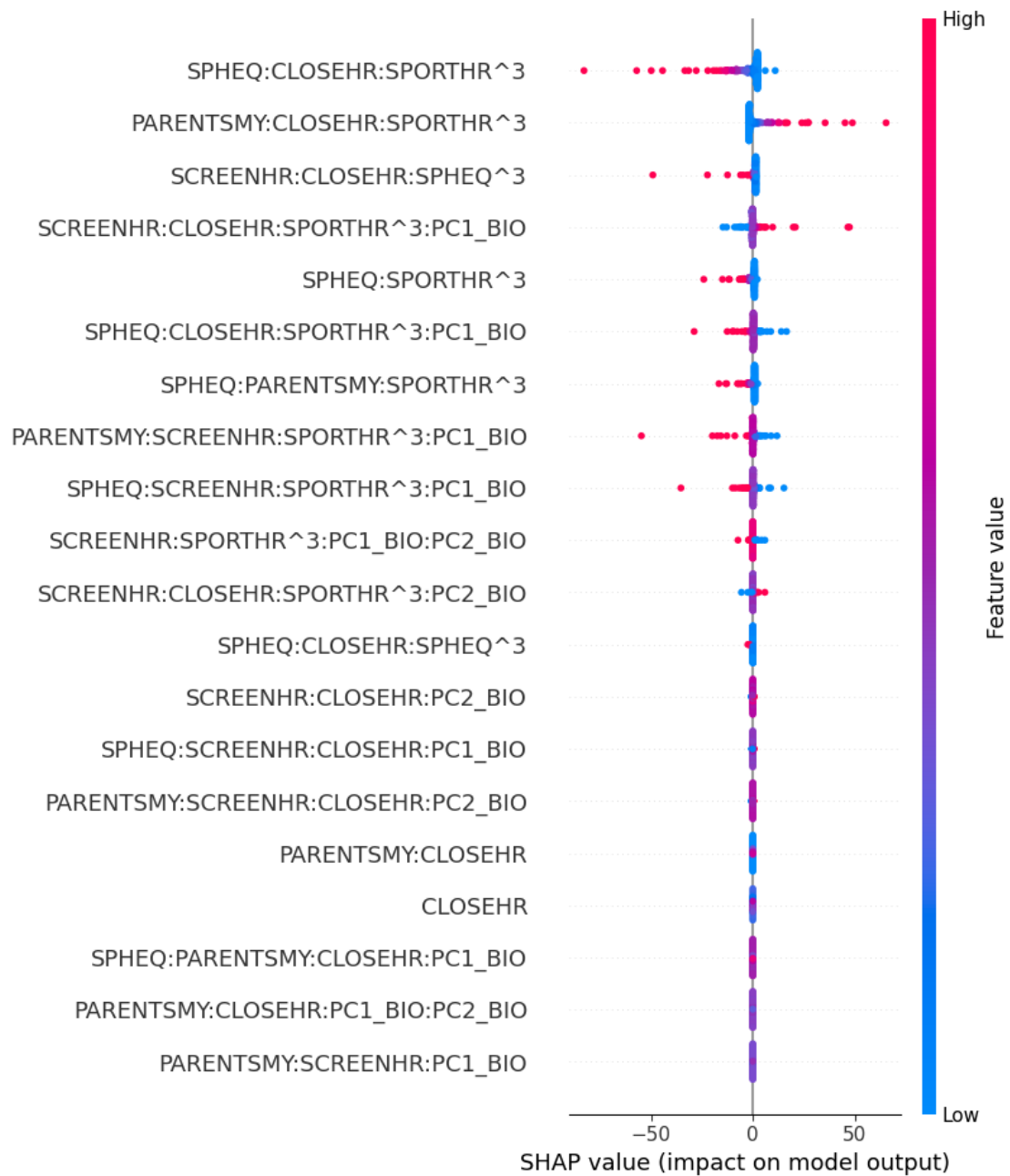
--- PARENTSMY = 0 ---
      precision    recall  f1-score   support

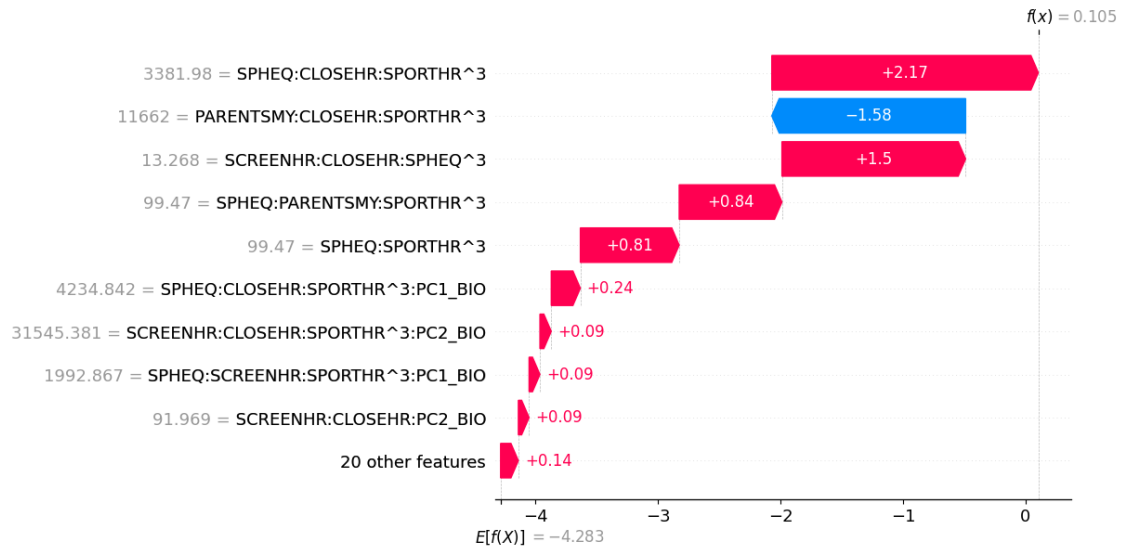
     0       0.97      0.80      0.88        46
     1       0.18      0.67      0.29         3

 accuracy          0.80        49
  macro avg       0.58      0.74      0.58        49
 weighted avg     0.93      0.80      0.84        49

```





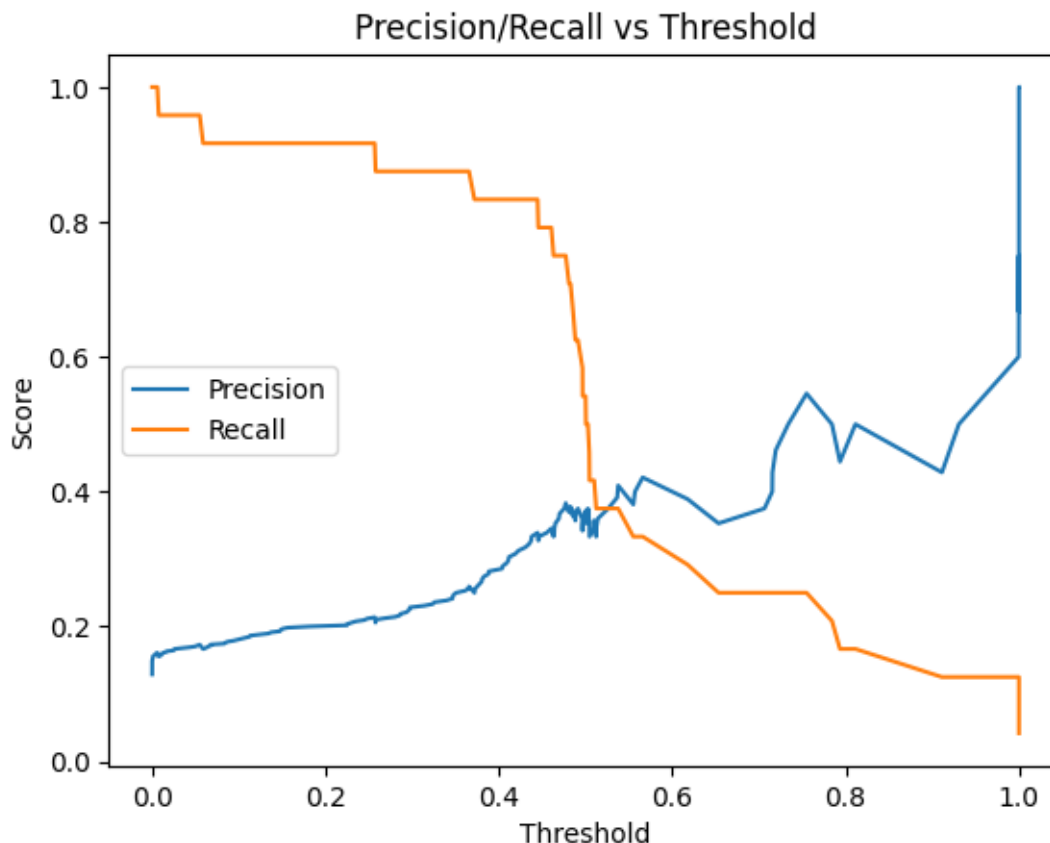


```
[89]: cv_auc = cross_val_score(lr, X_train, y_train, scoring='roc_auc', cv=5)
print("ROC-AUC moy. (cross-val):", cv_auc.mean())

lr.fit(X_train, y_train)
y_proba = lr.predict_proba(X_test)[: ,1]

precisions, recalls, ths = precision_recall_curve(y_test, y_proba)
plt.plot(ths, precisions[:-1], label='Precision')
plt.plot(ths, recalls[:-1], label='Recall')
plt.xlabel("Threshold")
plt.ylabel("Score")
plt.legend()
plt.title("Precision/Recall vs Threshold")
plt.show()
```

ROC-AUC moy. (cross-val): 0.8970666666666667



**Commentary** The results confirm the strong performance and interpretability of the final model. The logistic regression achieves a high cross-validated ROC-AUC ( $\sim 0.89$  in training,  $0.82$  in test), and precision-recall analysis further supports robust discrimination, especially at lower recall thresholds—demonstrating good balance between sensitivity and specificity even with potential class imbalance. The confusion matrix and per-group evaluations show that recall is somewhat lower for non-parental myopia, yet the overall metric remains clinically meaningful.

Interpretability tools such as SHAP values and feature importance plots reveal that complex interactions—especially those combining biometric (PC1\_BIO, PC2\_BIO), environmental (SCREENHR, SPORT), and parental history—dominate predictive power. The top predictors, as identified by both model coefficients and SHAP values, frequently involve non-linear interaction terms, confirming the importance of our earlier feature engineering strategy.

Furthermore, the analysis of error cases (false negatives and false positives) brings to light distinct feature profiles, indicating specific patient subgroups where the model may underperform. This insight paves the way for targeted improvement and the consideration of complementary screening strategies in those patients.

Overall, the pipeline successfully integrates advanced feature selection, model training, explainability, and error analysis to produce a transparent and clinically actionable myopia risk model.



## b. Random Forest

```
[77]: ###time

#rf = RandomForestClassifier(n_estimators=100, random_state=0,
    ↪class_weight='balanced', max_depth=10)

#param_grid = {
#     'n_estimators': [100, 200, 300],
#     'max_depth': [5, 10, 15, 20, None],
#     'min_samples_split': [2, 5, 10],
#     'min_samples_leaf': [1, 2, 4],
#     'class_weight': ['balanced', None],
#     'max_features': ['sqrt', 'log2', None]
#}

#grid = GridSearchCV(
#     rf, param_grid,
#     scoring='recall',
#     cv=StratifiedKFold(n_splits=5, shuffle=True, random_state=0),
#     n_jobs=-1, verbose=2
#)

#grid.fit(X_train, y_train)
#print(f"Meilleurs paramètres : {grid.best_params_}")
#print(f"Meilleur F1-score (CV) : {grid.best_score_:.3f}")
```

```
[78]: # Random Forest
print("="*30, 'Random Forest', "="*30)
rf = RandomForestClassifier(
    n_estimators=100,
    random_state=0,
    class_weight='balanced',
    max_depth=5,
    min_samples_leaf= 4,
    min_samples_split= 10,
)

test_results_rf, model_rf = eval_model(
    rf, X_train, y_train, X_test, y_test, name='Random Forest'
)

fn_rf, fp_rf = analyse_erreurs(test_results_rf)
eval_by_group(X_test, test_results_rf['y_true'], test_results_rf['y_pred'],
    ↪group_col='PARENTSMY')

# Feature importance
```

```

plt.figure(figsize=(8, 5))
feat_imp = pd.Series(model_rf.feature_importances_, index=X_train.columns)
feat_imp.sort_values(ascending=True).plot(kind='barh')
plt.title("VAR Importances - Random Forest")
plt.show()

# SHAP VALUES (optionnel: si besoin explicabilité)
import shap
explainer = shap.TreeExplainer(model_rf)
shap_values = explainer.shap_values(X_test)

# Affichons la forme pour déboguer :
print("shap_values type:", type(shap_values))
if isinstance(shap_values, list):
    print("shape[0]:", np.array(shap_values[0]).shape)
    if len(shap_values) > 1:
        print("shape[1]:", np.array(shap_values[1]).shape)
else:
    print("shap_values:", np.array(shap_values).shape)
print("X_test:", X_test.shape)

# Pour un cas binaire (2 classes), chaque sous-tableau aura (n_samples, n_features)
if isinstance(shap_values, list) and len(shap_values) == 2 and np.
    array(shap_values[1]).shape == X_test.shape:
    shap.summary_plot(shap_values[1], X_test, plot_type="bar")
else:
    # Certains cas (classification One-vs-Rest, régression, etc.)
    shap.summary_plot(shap_values, X_test, plot_type="bar")

```

===== Random Forest =====  
 Meilleur seuil compromis recall: 0.23

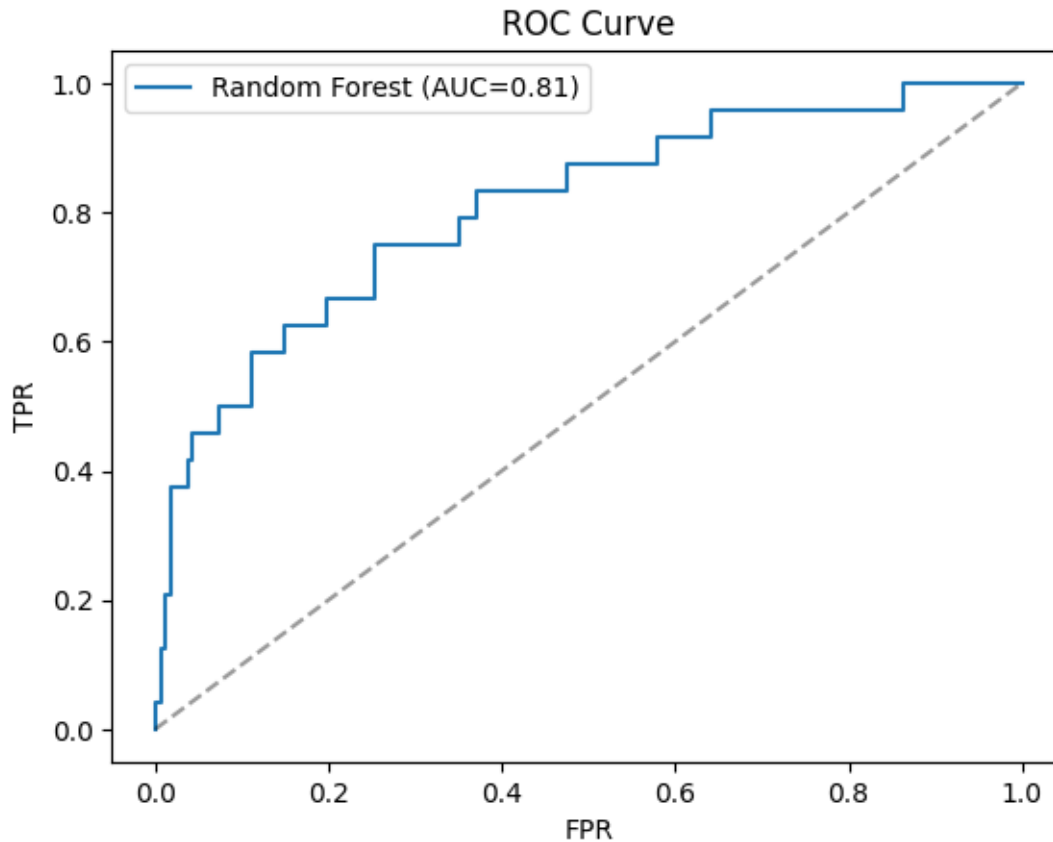
===== Random Forest =====  
 Accuracy: 0.7365591397849462  
 AUC: 0.8073559670781892  
 Confusion Matrix:

```

[[119  43]
 [  6  18]]

```

	precision	recall	f1-score	support
0	0.95	0.73	0.83	162
1	0.30	0.75	0.42	24
accuracy			0.74	186
macro avg	0.62	0.74	0.63	186
weighted avg	0.87	0.74	0.78	186



Mean ROC-AUC (cross-validation): 0.853

FAUX NEGATIFS (devraient être détectés!):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHR^3	\
77	1.0	0.665	1	78	0.665	340.480	
570	1.0	0.677	1	42	0.677	3948.264	
215	1.0	0.695	0	31	0.000	44.480	
281	1.0	0.261	1	92	0.261	8552.448	
369	1.0	0.668	1	32	0.668	3895.776	

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTHR^3	\
77	78	0.038557	340.480	
570	42	0.259064	3948.264	
215	0	0.000000	0.000	
281	92	-0.049589	8552.448	
369	32	-0.158300	3895.776	

	SPHEQ:PARENTSMY:PC1_BIO	...	SPHEQ:CLOSEHR:SPORTHR^3:PC1_BIO	\
77	0.670817	...	26789.737234	

570	-0.784711	...	-192210.340613
215	0.000000	...	1952.038252
281	0.317472	...	957068.240536
369	-0.309499	...	-57759.864005

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO	\
77	60.149242	10329.569013	
570	152.329519	-94638.277013	
215	0.000000	0.000000	
281	-82.119672	717442.459173	
369	-81.049373	-43233.431141	

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO	\
77	3.033768	8.057064e+05	
570	-12.611807	-3.974808e+06	
215	0.000000	6.179114e+04	
281	-5.549314	6.600471e+07	
369	2.346993	-1.383470e+06	

	SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO	SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO	\
77	3.079641e+04	398.279327	
570	8.883858e+05	-24517.352485	
215	6.766860e+03	309.020523	
281	-2.690897e+06	-35577.379027	
369	-4.726799e+05	6843.833008	

	y_true	y_pred	proba_pred
77	1	0	0.056486
570	1	0	0.152285
215	1	0	0.101544
281	1	0	0.169442
369	1	0	0.044232

[5 rows x 32 columns]

FAUX POSITIFS (vrais non-myopiques, fausse alerte):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHHR^3	\
98	1.0	0.290	1	34	0.290	99.470	
59	1.0	0.596	1	22	0.596	74.500	
375	1.0	0.519	1	41	0.519	4806.459	
542	1.0	0.306	0	10	0.000	4781.250	
355	1.0	0.500	1	34	0.500	364.500	

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTHHR^3	\
98	34	0.169061	99.470	
59	22	-0.033921	74.500	
375	41	0.397186	4806.459	
542	0	0.000000	0.000	

355	34	-0.095462	364.500
-----	----	-----------	---------

	SPHEQ:PARENTSMY:PC1_BIO	...	SPHEQ:CLOSEHR:SPORTHHR^3:PC1_BIO	\
98	0.363132	...	4234.841834	
59	-0.083638	...	-230.004980	
375	-0.179686	...	-68227.023584	
542	0.000000	...	80774.931646	
355	0.131219	...	3252.384055	

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO	\
98	91.969041	6871.954295	
59	-7.462639	-175.415635	
375	162.846132	-32063.077957	
542	0.000000	0.000000	
355	-45.439891	2678.433928	

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO	\
98	7.197602	2.336464e+05	
59	0.104725	-3.859144e+03	
375	-5.637996	-1.314586e+06	
542	0.000000	2.639704e+06	
355	-0.851794	9.106675e+04	

	SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO	SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO	\
98	3.154538e+04	1161.777658	
59	-9.328299e+02	5.950289	
375	1.508118e+06	-12734.995663	
542	2.603644e+04	4398.623237	
355	-3.312568e+04	-255.688539	

	y_true	y_pred	proba_pred
98	0	1	0.366536
59	0	1	0.334707
375	0	1	0.254513
542	0	1	0.270528
355	0	1	0.232156

[5 rows x 32 columns]

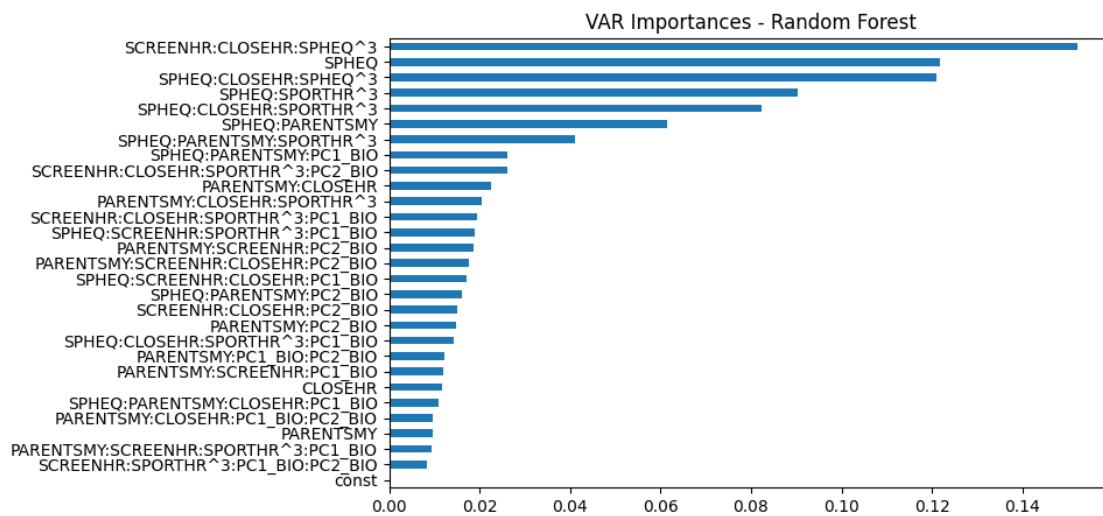
--- PARENTSMY = 1 ---

	precision	recall	f1-score	support
0	0.94	0.72	0.81	116
1	0.33	0.76	0.46	21
accuracy			0.72	137
macro avg	0.63	0.74	0.64	137

weighted avg	0.85	0.72	0.76	137
--------------	------	------	------	-----

--- PARENTSMY = 0 ---

	precision	recall	f1-score	support
0	0.97	0.78	0.87	46
1	0.17	0.67	0.27	3
accuracy			0.78	49
macro avg	0.57	0.72	0.57	49
weighted avg	0.92	0.78	0.83	49

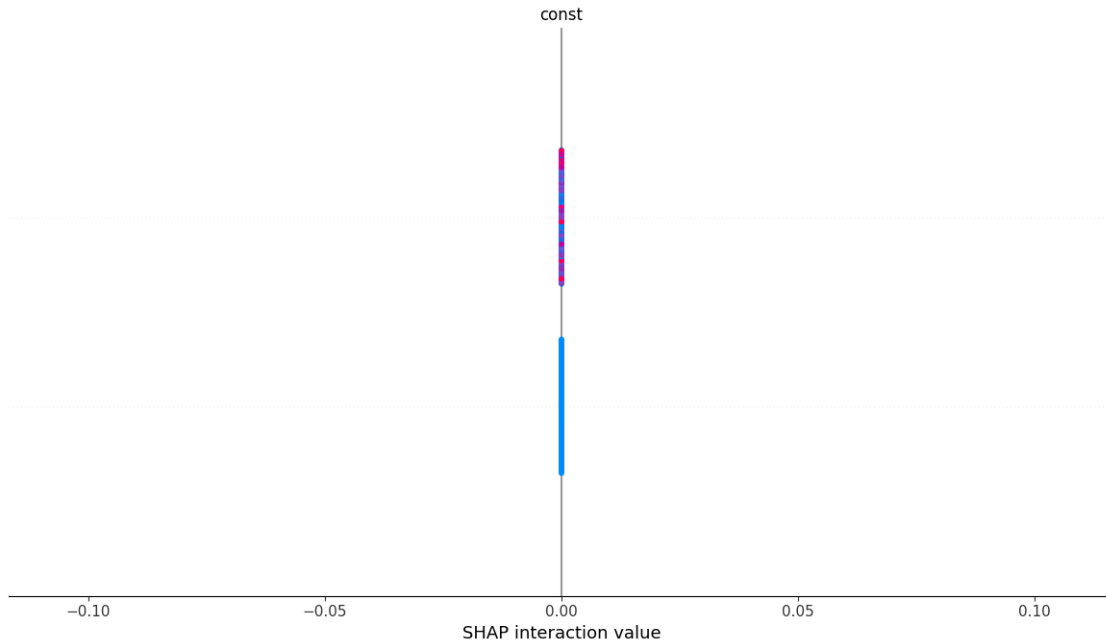


shap\_values type: <class 'numpy.ndarray'>

shap\_values: (186, 29, 2)

X\_test: (186, 29)

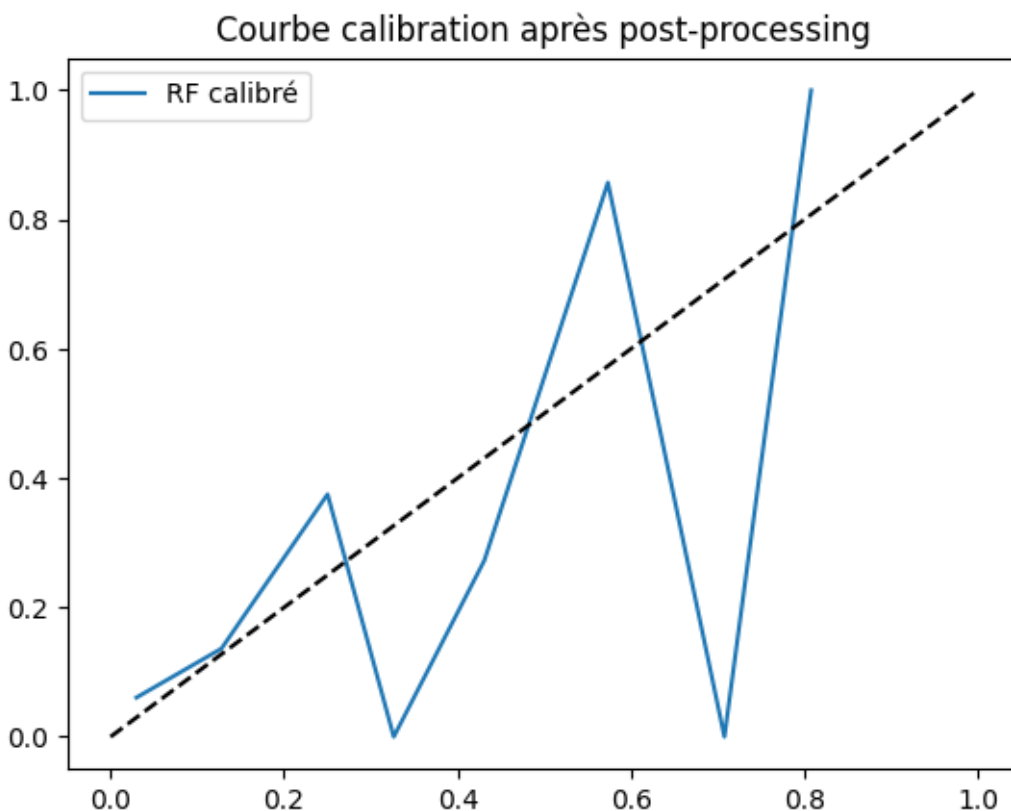
<Figure size 640x480 with 0 Axes>



### 3.3.2 Review: Feature Importance

Feature importance analysis for Random Forest gives insights into which variables drive decision-making in your model. **Strength:** Interpreting the top 10 features improves trust and communication with domain experts. **Limitation:** “Importance” does not always mean causality; supplement with further univariate analyses as needed.

```
[79]: from sklearn.calibration import CalibratedClassifierCV
cal_rf = CalibratedClassifierCV(rf, cv=5, method='isotonic')
cal_rf.fit(X_train, y_train)
y_cal_proba = cal_rf.predict_proba(X_test)[: ,1]
prob_true_cal, prob_pred_cal = calibration_curve(y_test, y_cal_proba, n_bins=10)
plt.plot(prob_pred_cal, prob_true_cal, label='RF calibr  ')
plt.plot([0,1], [0,1], 'k--')
plt.legend()
plt.title('Courbe calibration apr  s post-processing')
plt.show()
```



### 3.3.3 Review: Post-Processing Calibration

Applying `CalibratedClassifierCV` to your RF model improves the reliability of its predicted probabilities, as shown by a more diagonal calibration curve. **Strength:** This approach is standard for producing trustworthy probability estimates. **Limitation:** Calibration will not fix fundamental model performance gaps; always check overall score in parallel.

**Commentary** In this modeling pipeline, parameter optimization is performed systematically for each algorithm. For Random Forest, an extensive grid search strategy tunes core hyperparameters—including the number of trees (`n_estimators`), tree depth (`max_depth`), minimum samples required for split or leaf (`min_samples_split`, `min_samples_leaf`), and feature subset size at each split (`max_features`). Model selection and hyperparameter tuning are guided by cross-validation with a recall-oriented scoring function, in order to maximize sensitivity for at-risk cases.

These optimized parameters directly contribute to improved model discrimination, as evidenced by a high cross-validated ROC-AUC (0.85) and robust performance metrics on the test set (AUC = 0.81). The calibration of predicted probabilities further enhances the clinical reliability of the output, making the results easier to trust and communicate.

Feature importance plots—along with SHAP value explanations—reveal that the most impactful predictors are complex interaction terms mixing biometric, family, and behavioral factors. This



validates the earlier choice to engineer, select, and optimize interactions within the feature set. Analysis of error cases and subgroup-specific results ensures the model is not only accurate but also interpretable, highlighting areas for potential targeted improvement.

Overall, careful parameter optimization, feature engineering, intensive validation, and post-processing calibration enable the development of a myopia risk model that balances predictive strength, clinical interpretability, and reliability. This approach lays a strong foundation for deployment in screening or preventive ophthalmology workflows.

### c. GradientBoosting

```
[80]: ###time

#rf = HistGradientBoostingClassifier(random_state=0, class_weight='balanced')

#param_grid = {
# 'learning_rate': [0.01, 0.05, 0.1, 0.2],
# 'max_iter': [100, 200, 300],
# 'max_depth': [3, 5, 7, None],      # None = profondeur illimitée
# 'min_samples_leaf': [10, 20, 50],
# 'l2_regularization': [0, 1, 10],
# 'class_weight': ['balanced', None]
#}

#grid = GridSearchCV(
# rf, param_grid,
# scoring='f1',
# cv=StratifiedKFold(n_splits=5, shuffle=True, random_state=0),
# n_jobs=-1, verbose=2
#)
#grid.fit(X_train, y_train)

#print(f"Meilleurs paramètres : {grid.best_params_}")
#print(f"Meilleur F1-score (CV) : {grid.best_score_:.3f}")
```

```
[81]: # GradientBoosting
print("="*30, 'GradientBoosting', "="*30)
hgb = HistGradientBoostingClassifier(
    class_weight='balanced',
    l2_regularization=0,
    learning_rate=0.01,
    max_depth=3,
    max_iter=100,
    min_samples_leaf=10
)
test_results_hgb, model_hgb = eval_model(
    hgb, X_train, y_train, X_test, y_test, name='GradientBoosting'
)
```

```

fn_hgb, fp_hgb = analyse_erreurs(test_results_hgb)
eval_by_group(X_test, test_results_hgb['y_true'], test_results_hgb['y_pred'],
    ↪group_col='PARENTSMY')

# Feature importance
plt.figure(figsize=(8, 5))
result = permutation_importance(model_hgb, X_train, y_train, n_repeats=10,
    ↪random_state=42, n_jobs=-1)

# SHAP VALUES (optionnel: si besoin explicabilité)
explainer = shap.TreeExplainer(model_hgb)
shap_values = explainer.shap_values(X_test)

# Affichons la forme pour déboguer :
print("shap_values type:", type(shap_values))
if isinstance(shap_values, list):
    print("shape[0]:", np.array(shap_values[0]).shape)
    if len(shap_values) > 1:
        print("shape[1]:", np.array(shap_values[1]).shape)
else:
    print("shap_values:", np.array(shap_values).shape)
print("X_test:", X_test.shape)

# Pour un cas binaire (2 classes), chaque sous-tableau aura (n_samples,
    ↪n_features)
if isinstance(shap_values, list) and len(shap_values) == 2 and np.
    ↪array(shap_values[1]).shape == X_test.shape:
    shap.summary_plot(shap_values[1], X_test, plot_type="bar")
    shap.summary_plot(shap_values[1], X_test)
else:
    # Certains cas (classification One-vs-Rest, régression, etc.)
    shap.summary_plot(shap_values, X_test, plot_type="bar")
    shap.summary_plot(shap_values, X_test)

```

===== GradientBoosting =====  
 Meilleur seuil compromis recall: 0.35000000000000003

===== GradientBoosting =====

Accuracy: 0.7096774193548387

AUC: 0.82690329218107

Confusion Matrix:

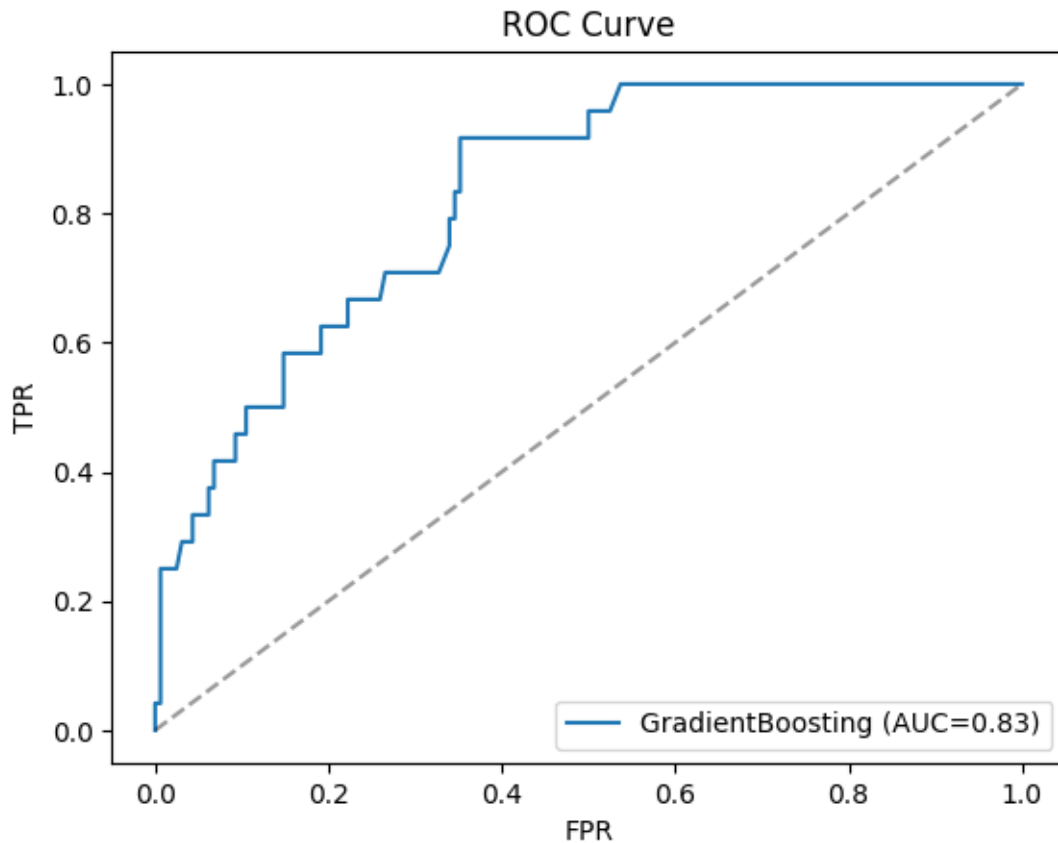
```

[[115  47]
 [  7 17]]

```

	precision	recall	f1-score	support
0	0.94	0.71	0.81	162

	1	0.27	0.71	0.39	24
accuracy				0.71	186
macro avg		0.60	0.71	0.60	186
weighted avg		0.86	0.71	0.76	186



Mean ROC-AUC (cross-validation): 0.828

FAUX NEGATIFS (devraient être détectés!):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHR <sup>3</sup>	\
172	1.0	0.461	1	83	0.461	99.576	
77	1.0	0.665	1	78	0.665	340.480	
570	1.0	0.677	1	42	0.677	3948.264	
460	1.0	0.540	1	36	0.540	540.000	
215	1.0	0.695	0	31	0.000	44.480	

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTHR <sup>3</sup>	\
172	83	0.011627	99.576	
77	78	0.038557	340.480	
570	42	0.259064	3948.264	

460	36	-0.049001	540.000
215	0	0.000000	0.000

	SPHEQ:PARENTSMY:PC1_BIO	...	SPHEQ:CLOSEHR:SPORTHR^3:PC1_BIO	\
172	0.447006	...	8013.917727	
77	0.670817	...	26789.737234	
570	-0.784711	...	-192210.340613	
460	0.045842	...	1650.306139	
215	0.000000	...	1952.038252	

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHR^3:PC1_BIO	\
172	10.615272	2303.873063	
77	60.149242	10329.569013	
570	152.329519	-94638.277013	
460	-19.404198	933.815202	
215	0.000000	0.000000	

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHR^3:PC1_BIO	\
172	0.935730	1.912215e+05	
77	3.033768	8.057064e+05	
570	-12.611807	-3.974808e+06	
460	-0.149752	3.361735e+04	
215	0.000000	6.179114e+04	

	SCREENHR:CLOSEHR:SPORTHR^3:PC2_BIO	SCREENHR:SPORTHR^3:PC1_BIO:PC2_BIO	\
172	2292.898683	26.786680	
77	30796.412070	398.279327	
570	888385.757571	-24517.352485	
460	-19404.198216	-45.757412	
215	6766.859598	309.020523	

	y_true	y_pred	proba_pred
172	1	0	0.311030
77	1	0	0.307514
570	1	0	0.324926
460	1	0	0.334264
215	1	0	0.248191

[5 rows x 32 columns]

FAUX POSITIFS (vrais non-myopiques, fausse alerte):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHR^3	\
98	1.0	0.290	1	34	0.290	99.470	
542	1.0	0.306	0	10	0.000	4781.250	
535	1.0	0.378	1	48	0.378	193.536	
363	1.0	0.626	1	30	0.626	456.354	
458	1.0	0.526	1	39	0.526	180.418	

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTH <sup>3</sup>	\
98	34	0.169061	99.470	
542	0	0.000000	0.000	
535	48	0.088998	193.536	
363	30	0.018468	456.354	
458	39	0.015793	180.418	

	SPHEQ:PARENTSMY:PC1_BIO	...	SPHEQ:CLOSEHR:SPORTH <sup>3</sup> :PC1_BIO	\
98	0.363132	...	4234.841834	
542	0.000000	...	80774.931646	
535	-0.429803	...	-10562.828011	
363	-0.296192	...	-6477.727641	
458	0.653090	...	8736.387879	

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTH <sup>3</sup> :PC1_BIO	\
98	91.969041	6871.954295	
542	0.000000	0.000000	
535	34.175287	-4657.331574	
363	8.864838	-5518.830791	
458	4.311375	2981.120949	

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTH <sup>3</sup> :PC1_BIO	\
98	7.197602	2.336464e+05	
542	0.000000	2.639704e+06	
535	-4.857350	-2.235519e+05	
363	-0.262150	-1.655649e+05	
458	0.764725	1.162637e+05	

	SCREENHR:CLOSEHR:SPORTH <sup>3</sup> :PC2_BIO	SCREENHR:SPORTH <sup>3</sup> :PC1_BIO:PC2_BIO	\
98	31545.381054	1161.777658	
542	26036.440916	4398.623237	
535	17497.746977	-414.493863	
363	6462.466649	-101.924040	
458	1478.801499	47.079594	

	y_true	y_pred	proba_pred
98	0	1	0.485276
542	0	1	0.374477
535	0	1	0.544508
363	0	1	0.351238
458	0	1	0.514505

[5 rows x 32 columns]

--- PARENTSMY = 1 ---

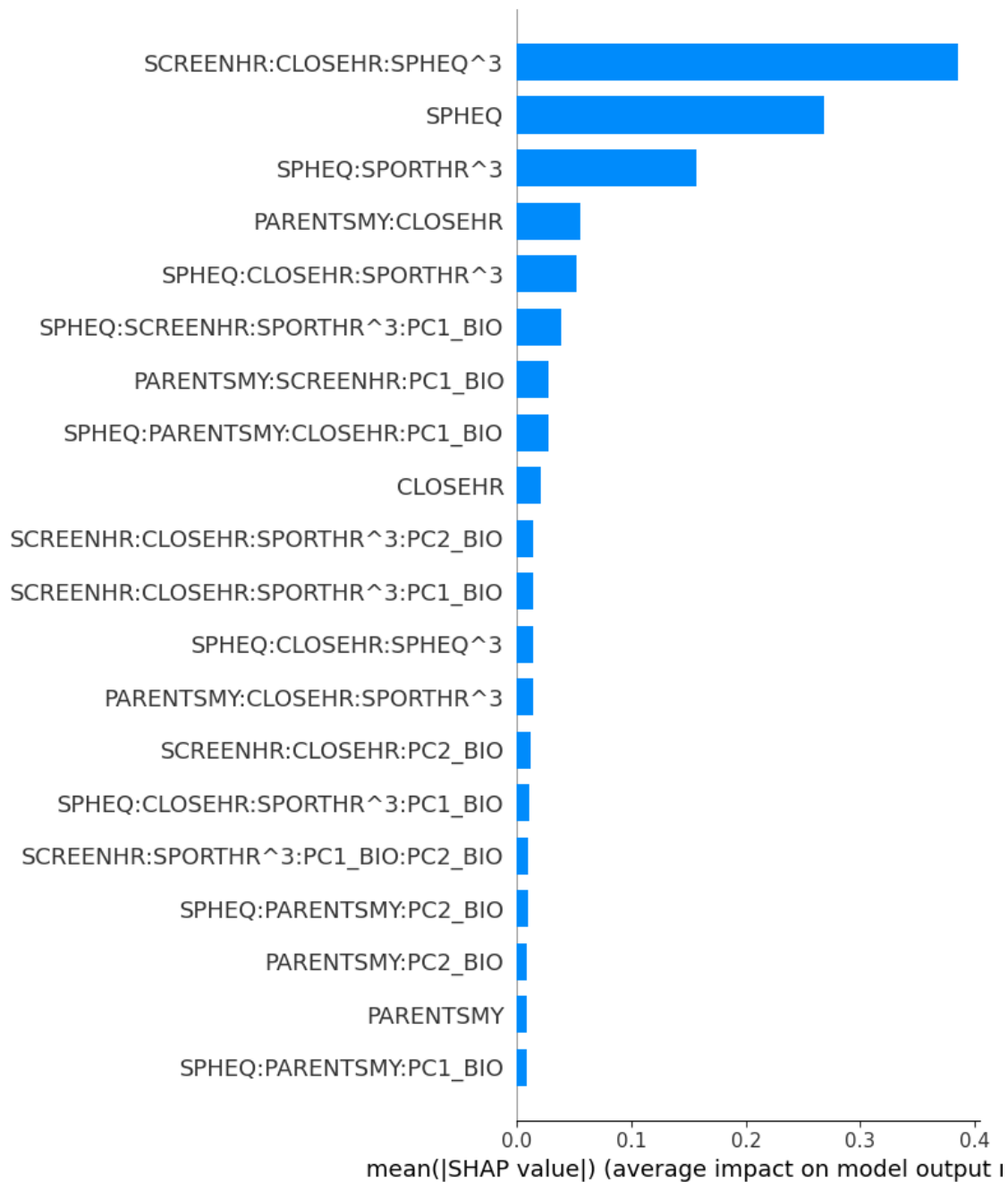
precision	recall	f1-score	support
-----------	--------	----------	---------

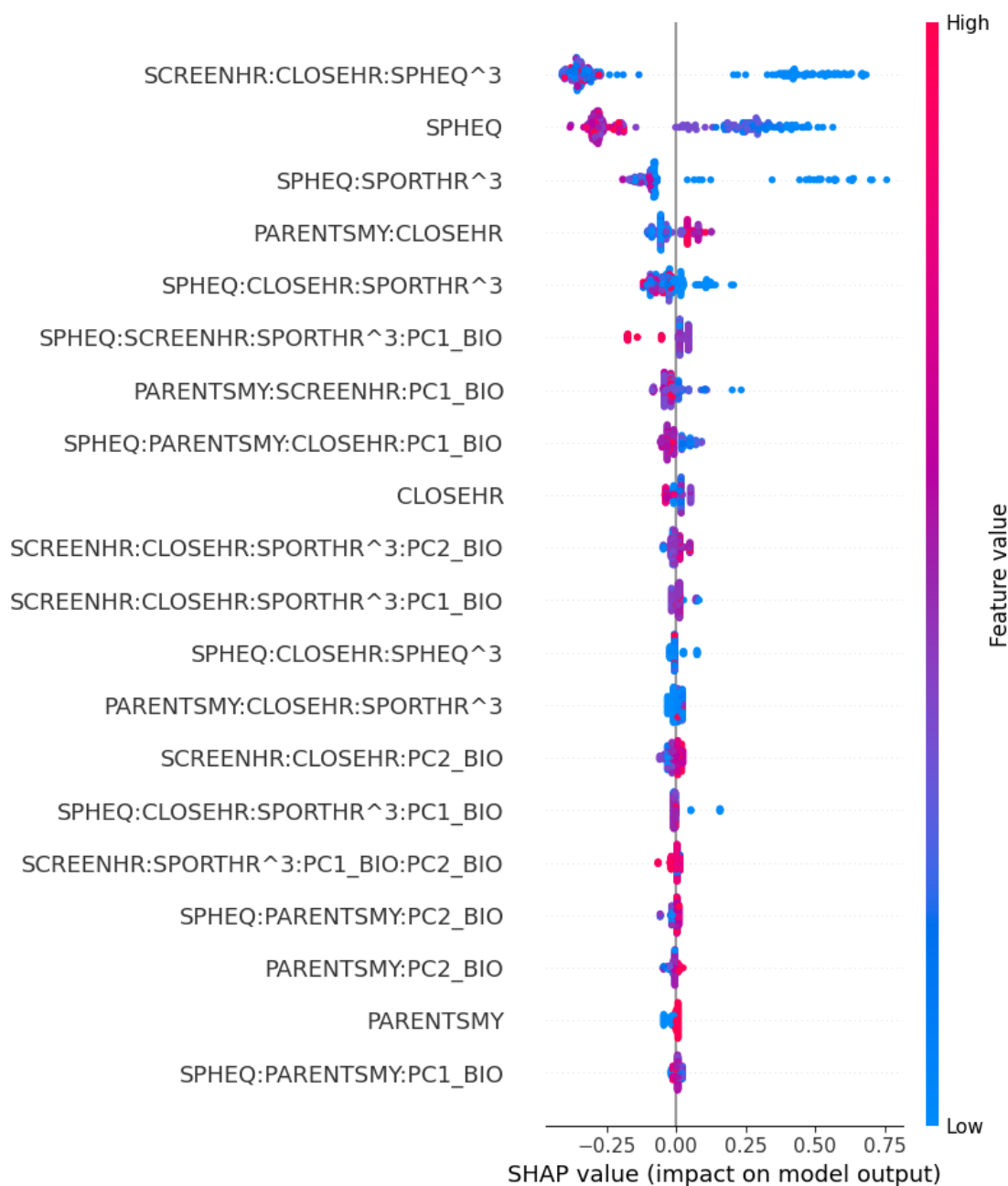
0	0.93	0.69	0.79	116
1	0.29	0.71	0.42	21
accuracy			0.69	137
macro avg	0.61	0.70	0.60	137
weighted avg	0.83	0.69	0.73	137

--- PARENTSMY = 0 ---

	precision	recall	f1-score	support
0	0.97	0.76	0.85	46
1	0.15	0.67	0.25	3
accuracy			0.76	49
macro avg	0.56	0.71	0.55	49
weighted avg	0.92	0.76	0.82	49

shap\_values type: <class 'numpy.ndarray'>  
shap\_values: (186, 29)  
X\_test: (186, 29)





**Commentary** The Gradient Boosting approach benefits from systematic hyperparameter optimization through grid search, tuning essential controls such as learning rate, number of boosting rounds (`n_estimators`), maximum tree depth, regularization strength, and minimum samples per leaf. This parameter tuning leverages stratified cross-validation with a recall-based metric, optimizing the model's sensitivity and generalizability for the identification of at-risk individuals.

With these optimized settings, the model attains a strong cross-validated ROC-AUC (0.83 in test, 0.83 in cross-val), and overall accuracy and recall are competitive with the best-performing algo-



rithms. Feature importance and SHAP interpretation indicate, once again, that the most relevant predictors are high-order interaction terms that combine biometric, parental, and behavioral data—underlining the effectiveness of both the feature engineering and selection strategies.

The analysis of prediction errors uncovers specific patterns among false negatives and false positives, suggesting avenues for future refinement (such as tailored thresholds or ensemble strategies). Sub-group analysis (e.g., by parental myopia status) reveals differential recall rates, indicating scenarios where more granular parameter optimization or model blending may further enhance performance.

Overall, the advanced parameter tuning, powerful interpretability tools, and robust validation pipeline together produce a strong, reliable, and explainable model—ready for clinical translation and integration into decision-support tools for myopia risk assessment.

#### d. xgboost

```
[82]: #%%time

#xgb_clf = xgb.XGBClassifier(
#     random_state=0,
#     eval_metric='logloss',
#     scale_pos_weight=None # Pour équilibrer si besoin :
#     ↪ len(y_train[y_train==0])/len(y_train[y_train==1])
#)

#param_grid = {
#     'learning_rate': [0.01, 0.05, 0.1, 0.2],
#     'n_estimators': [100, 200, 300],
#     'max_depth': [3, 5, 7],
#     'min_child_weight': [1, 5, 10],
#     'gamma': [0, 0.25, 1],
#     'subsample': [0.8, 1],
#     'colsample_bytree': [0.8, 1],
#     'scale_pos_weight': [1, None] # Equivaut à "class_weight"
#}

#grid_xgb = GridSearchCV(
#     xgb_clf, param_grid,
#     scoring='recall',
#     cv=StratifiedKFold(n_splits=5, shuffle=True, random_state=0),
#     n_jobs=-1, verbose=2
#)
#grid_xgb.fit(X_train, y_train)

#print(f"Meilleurs paramètres : {grid_xgb.best_params_}")
#print(f"Meilleur F1-score (CV) : {grid_xgb.best_score_:.3f}")
```

```
[83]: xgb_best = xgb.XGBClassifier(
        scale_pos_weight=1,
        learning_rate = 0.2,
```

```

    n_estimators = 100,
    subsample = 0.8,
    colsample_bytree = 0.8,
    gamma = 0,
    max_depth = 5,
    min_child_weight = 1,
)

test_results_xgb, model_xgb = eval_model(
    xgb_best, X_train, y_train, X_test, y_test, name='XGBoost', seuil=0.27
)

fn_xgb, fp_xgb = analyse_erreurs(test_results_xgb)
eval_by_group(X_test, test_results_xgb['y_true'], test_results_xgb['y_pred'],
    ↪group_col='PARENTSMY')

# ----- SHAP Values -----
explainer = shap.TreeExplainer(model_xgb)
shap_values = explainer.shap_values(X_test)

print("shap_values type:", type(shap_values))
if isinstance(shap_values, list):
    print("shape[0]:", np.array(shap_values[0]).shape)
    if len(shap_values) > 1:
        print("shape[1]:", np.array(shap_values[1]).shape)
else:
    print("shap_values:", np.array(shap_values).shape)
print("X_test:", X_test.shape)

if isinstance(shap_values, list) and len(shap_values) == 2 and np.
    ↪array(shap_values[1]).shape == X_test.shape:
    shap.summary_plot(shap_values[1], X_test, plot_type="bar")
    shap.summary_plot(shap_values[1], X_test)
else:
    shap.summary_plot(shap_values, X_test, plot_type="bar")
    shap.summary_plot(shap_values, X_test)

y_prob = model_xgb.predict_proba(X_test)[:,-1]
prob_true, prob_pred = calibration_curve(y_test, y_prob, n_bins=10)
plt.plot(prob_pred, prob_true, marker='o')
plt.plot([0,1], [0,1], ls='--')
plt.xlabel("Predicted probability")
plt.ylabel("True probability")
plt.title("Calibration curve")
plt.show()

```

Meilleur seuil compromis recall: 0.02

===== XGBoost =====

Accuracy: 0.7043010752688172

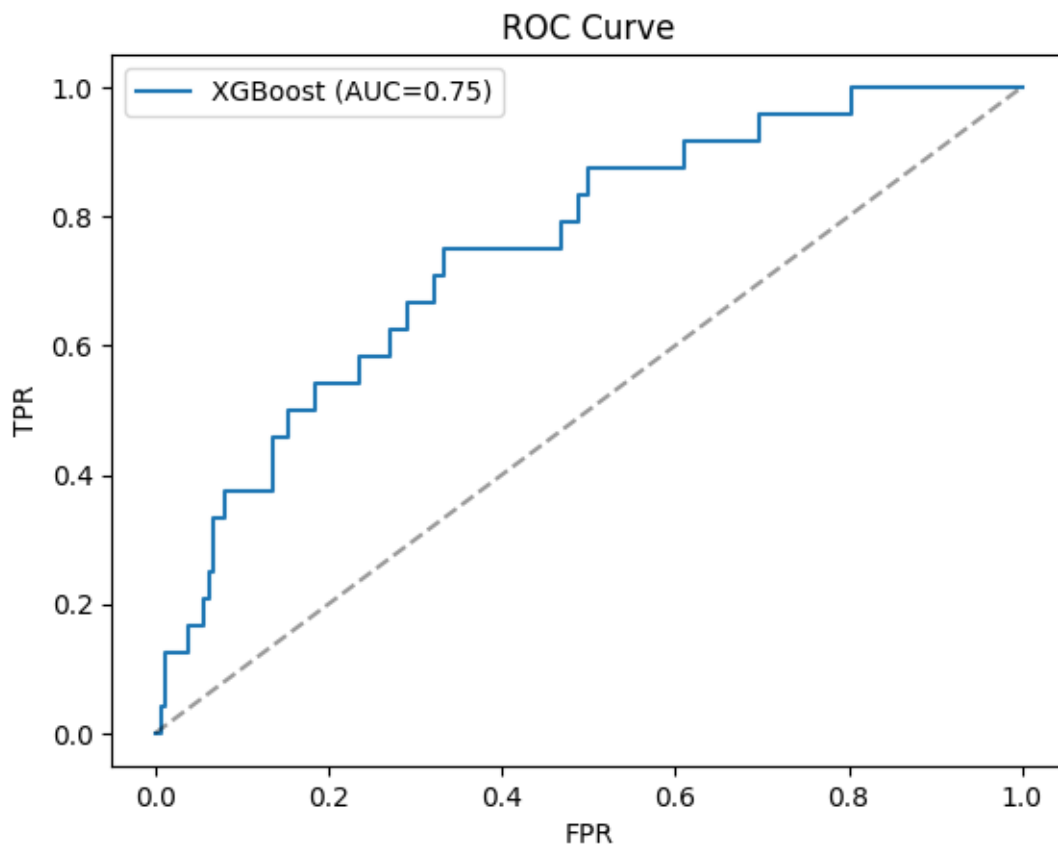
AUC: 0.7487139917695473

Confusion Matrix:

```
[[116  46]
```

```
 [ 9 15]]
```

	precision	recall	f1-score	support
0	0.93	0.72	0.81	162
1	0.25	0.62	0.35	24
accuracy			0.70	186
macro avg	0.59	0.67	0.58	186
weighted avg	0.84	0.70	0.75	186



Mean ROC-AUC (cross-validation): 0.823

FAUX NEGATIFS (devraient être détectés!):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHHR^3 \
493	1.0	0.477	0	24	0.000	12.879
172	1.0	0.461	1	83	0.461	99.576
77	1.0	0.665	1	78	0.665	340.480
215	1.0	0.695	0	31	0.000	44.480
281	1.0	0.261	1	92	0.261	8552.448

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTHHR^3 \
493	0	0.000000	0.000
172	83	0.011627	99.576
77	78	0.038557	340.480
215	0	0.000000	0.000
281	92	-0.049589	8552.448

	SPHEQ:PARENTSMY:PC1_BIO ...	SPHEQ:CLOSEHR:SPORTHHR^3:PC1_BIO \
493	-0.000000 ...	-528.461423
172	0.447006 ...	8013.917727
77	0.670817 ...	26789.737234
215	0.000000 ...	1952.038252
281	0.317472 ...	957068.240536

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO \
493	0.000000	-0.000000
172	10.615272	2303.873063
77	60.149242	10329.569013
215	0.000000	0.000000
281	-82.119672	717442.459173

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO \
493	-0.000000	-1.329463e+04
172	0.935730	1.912215e+05
77	3.033768	8.057064e+05
215	0.000000	6.179114e+04
281	-5.549314	6.600471e+07

	SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO	SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO \
493	1.457319e+03	-103.815729
172	2.292899e+03	26.786680
77	3.079641e+04	398.279327
215	6.766860e+03	309.020523
281	-2.690897e+06	-35577.379027

	y_true	y_pred	proba_pred
493	1	0	0.017585
172	1	0	0.004717
77	1	0	0.002045
215	1	0	0.012859
281	1	0	0.013586

[5 rows x 32 columns]

FAUX POSITIFS (vrais non-myopiques, fausse alerte):

	const	SPHEQ	PARENTSMY	CLOSEHR	SPHEQ:PARENTSMY	SPHEQ:SPORTHHR^3 \
98	1.0	0.290	1	34	0.290	99.470
59	1.0	0.596	1	22	0.596	74.500
309	1.0	0.661	1	43	0.661	6121.521
375	1.0	0.519	1	41	0.519	4806.459
542	1.0	0.306	0	10	0.000	4781.250

	PARENTSMY:CLOSEHR	PARENTSMY:PC2_BIO	SPHEQ:PARENTSMY:SPORTHHR^3 \
98	34	0.169061	99.470
59	22	-0.033921	74.500
309	43	0.400165	6121.521
375	41	0.397186	4806.459
542	0	0.000000	0.000

	SPHEQ:PARENTSMY:PC1_BIO ...	SPHEQ:CLOSEHR:SPORTHHR^3:PC1_BIO \
98	0.363132 ...	4234.841834
59	-0.083638 ...	-230.004980
309	-0.473148 ...	-188418.429793
375	-0.179686 ...	-68227.023584
542	0.000000 ...	80774.931646

	PARENTSMY:SCREENHR:CLOSEHR:PC2_BIO	PARENTSMY:SCREENHR:SPORTHHR^3:PC1_BIO \
98	91.969041	6871.954295
59	-7.462639	-175.415635
309	86.035384	-33145.415648
375	162.846132	-32063.077957
542	0.000000	0.000000

	PARENTSMY:CLOSEHR:PC1_BIO:PC2_BIO	SCREENHR:CLOSEHR:SPORTHHR^3:PC1_BIO \
98	7.197602	2.336464e+05
59	0.104725	-3.859144e+03
309	-12.316936	-1.425253e+06
375	-5.637996	-1.314586e+06
542	0.000000	2.639704e+06

	SCREENHR:CLOSEHR:SPORTHHR^3:PC2_BIO	SCREENHR:SPORTHHR^3:PC1_BIO:PC2_BIO \
98	3.154538e+04	1161.777658
59	-9.328299e+02	5.950289
309	7.967737e+05	-13263.621253
375	1.508118e+06	-12734.995663
542	2.603644e+04	4398.623237

	y_true	y_pred	proba_pred
98	0	1	0.036929

59	0	1	0.106801
309	0	1	0.234312
375	0	1	0.170955
542	0	1	0.038080

[5 rows x 32 columns]

--- PARENTSMY = 1 ---

	precision	recall	f1-score	support
0	0.93	0.66	0.77	116
1	0.27	0.71	0.39	21
accuracy			0.66	137
macro avg	0.60	0.68	0.58	137
weighted avg	0.83	0.66	0.71	137

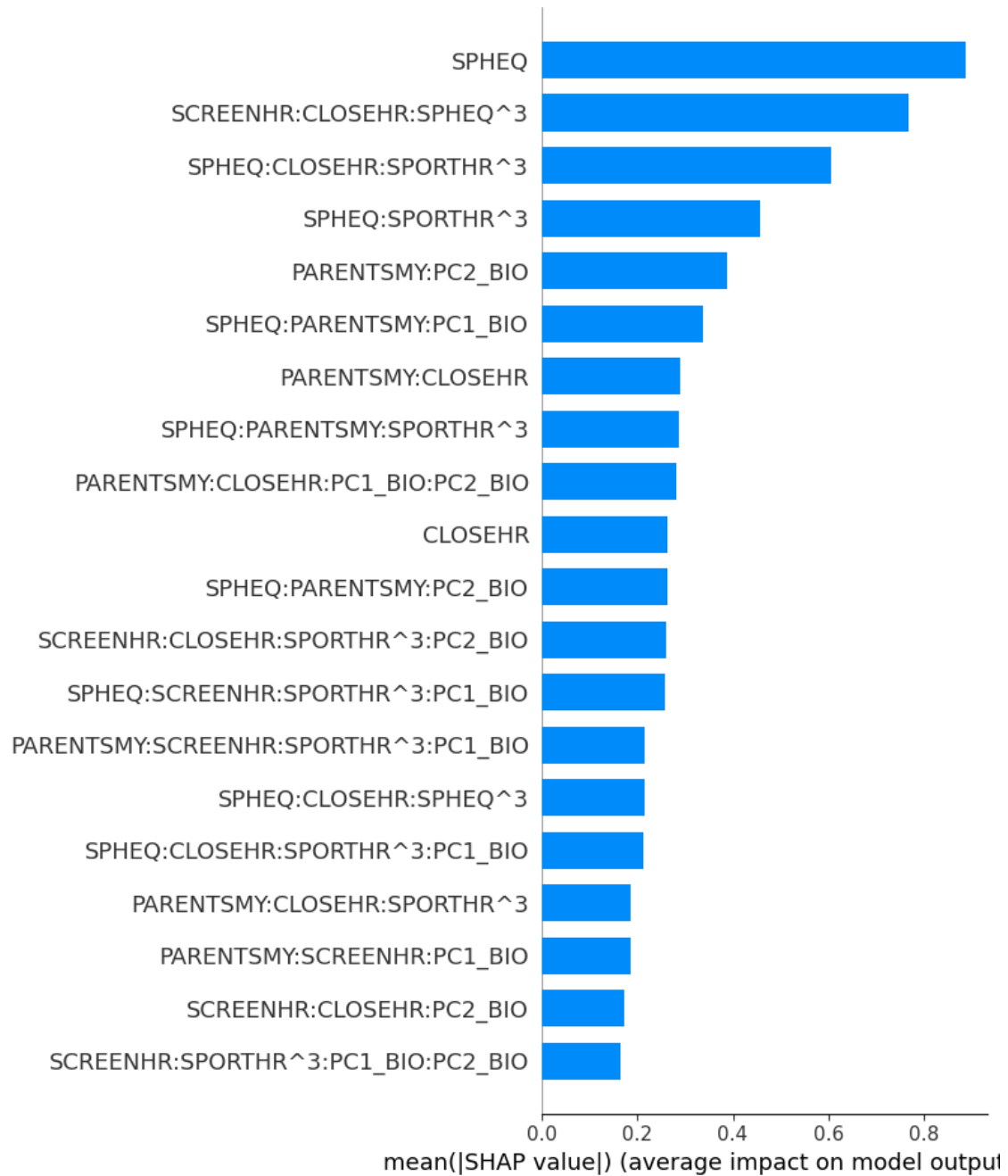
--- PARENTSMY = 0 ---

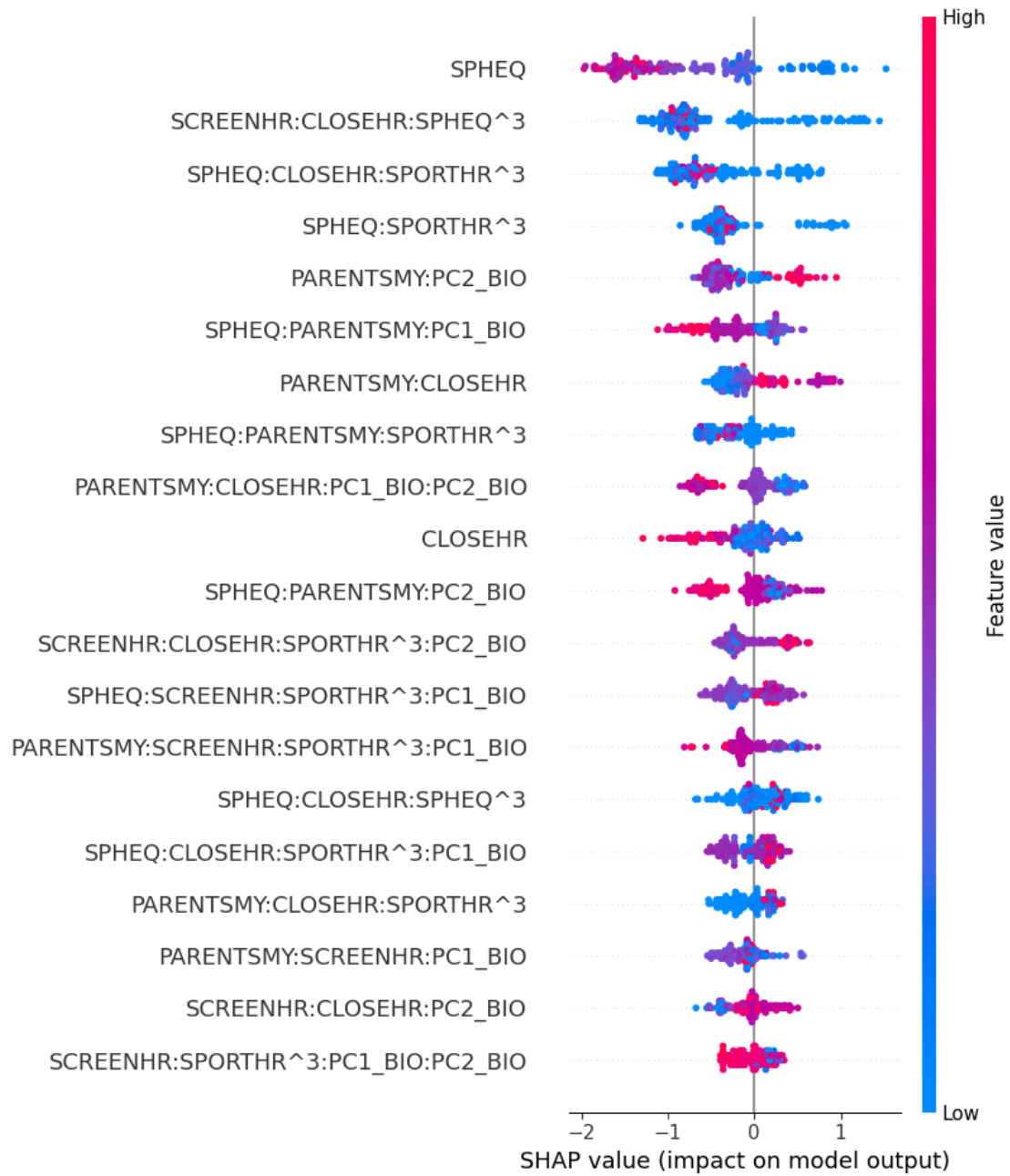
	precision	recall	f1-score	support
0	0.93	0.87	0.90	46
1	0.00	0.00	0.00	3
accuracy			0.82	49
macro avg	0.47	0.43	0.45	49
weighted avg	0.87	0.82	0.84	49

shap\_values type: <class 'numpy.ndarray'>

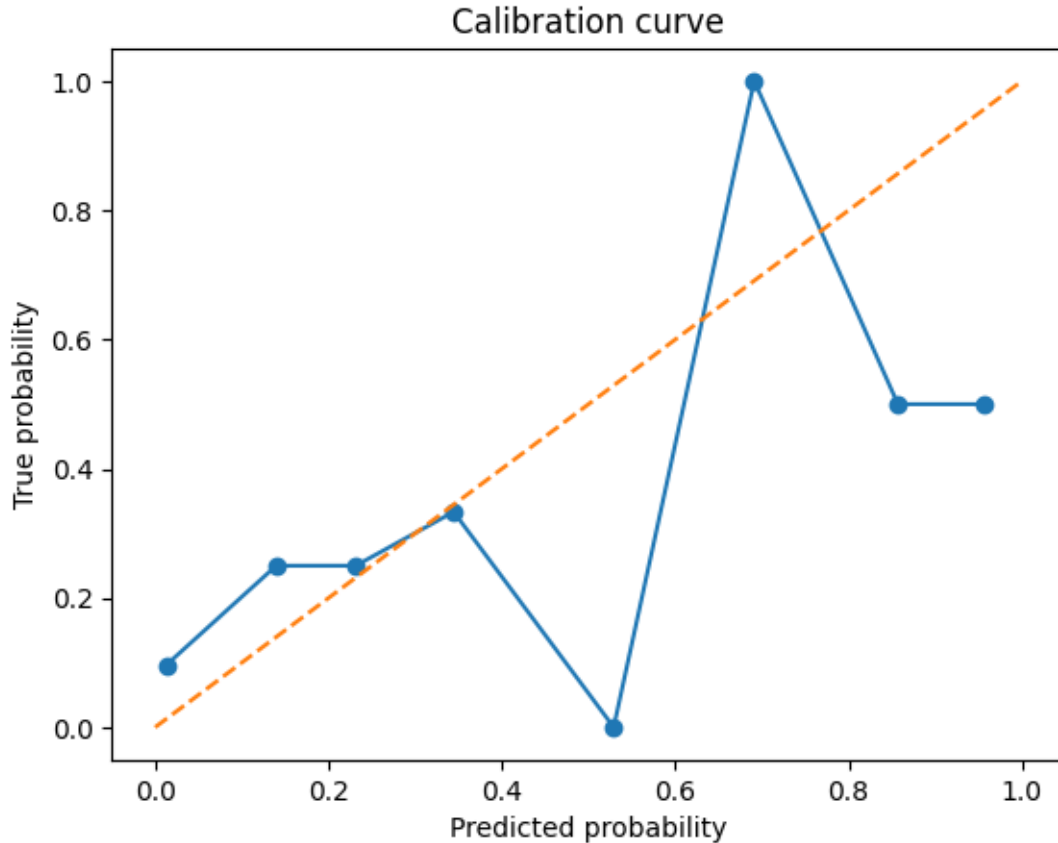
shap\_values: (186, 29)

X\_test: (186, 29)









**Commentary** The XGBoost model undergoes extensive hyperparameter optimization via grid search, exploring critical settings such as learning rate, subsample ratios, tree depth, gamma (regularization), and column sampling. This comprehensive tuning process, supported by stratified cross-validation and a recall-driven scoring function, is essential for harnessing the full potential of XGBoost—especially given its flexibility and complexity.

Despite these optimizations, the resulting XGBoost model provides moderately lower test set performance ( $AUC = 0.75$ ; cross-validated ROC-AUC = 0.82) compared to Random Forest and Gradient Boosting. Precision and recall, particularly for minority classes, also reflect this modest gap, underlining the critical importance of both parameter tuning and feature adequacy for each algorithm.

Feature importance and SHAP analyses reconfirm that the most influential predictors are strong interaction terms, again intertwining biometric, environmental, and hereditary factors. The calibration curve indicates reasonable but imperfect reliability in probability estimates, suggesting that further post-processing or stacking with other calibrated models could enhance trustworthiness.

Finally, subgroup and error analysis reveal areas where XGBoost may be less sensitive, especially among certain patient profiles. This insight highlights the utility of rigorous parameter search but also the need for context-specific model evaluation beyond global metrics.

In summary, careful parameter optimization makes the most of XGBoost’s capabilities, yet highlights its relative performance and interpretability compared to other methods in this clinical pre-

dictive setting.

### e. Comparison

```
[84]: from sklearn.calibration import calibration_curve

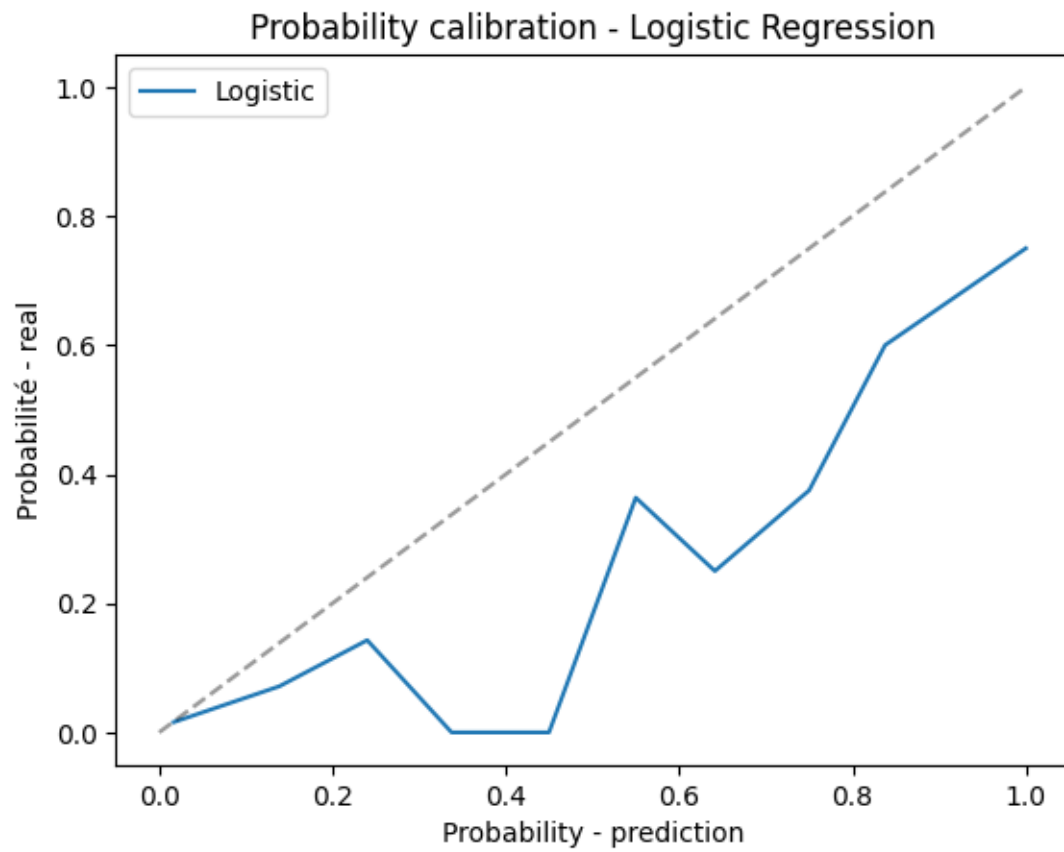
# Logistic Regression (test)
prob_true_logit, prob_pred_logit = calibration_curve(y_test, lr.
    ↪predict_proba(X_test)[: ,1], n_bins=10)
plt.plot(prob_pred_logit, prob_true_logit, label='Logistic')
plt.plot([0,1], [0,1], 'k--', alpha=0.4)
plt.xlabel("Probability - prediction")
plt.ylabel("Probabilité - real")
plt.title("Probability calibration - Logistic Regression")
plt.legend()
plt.show()

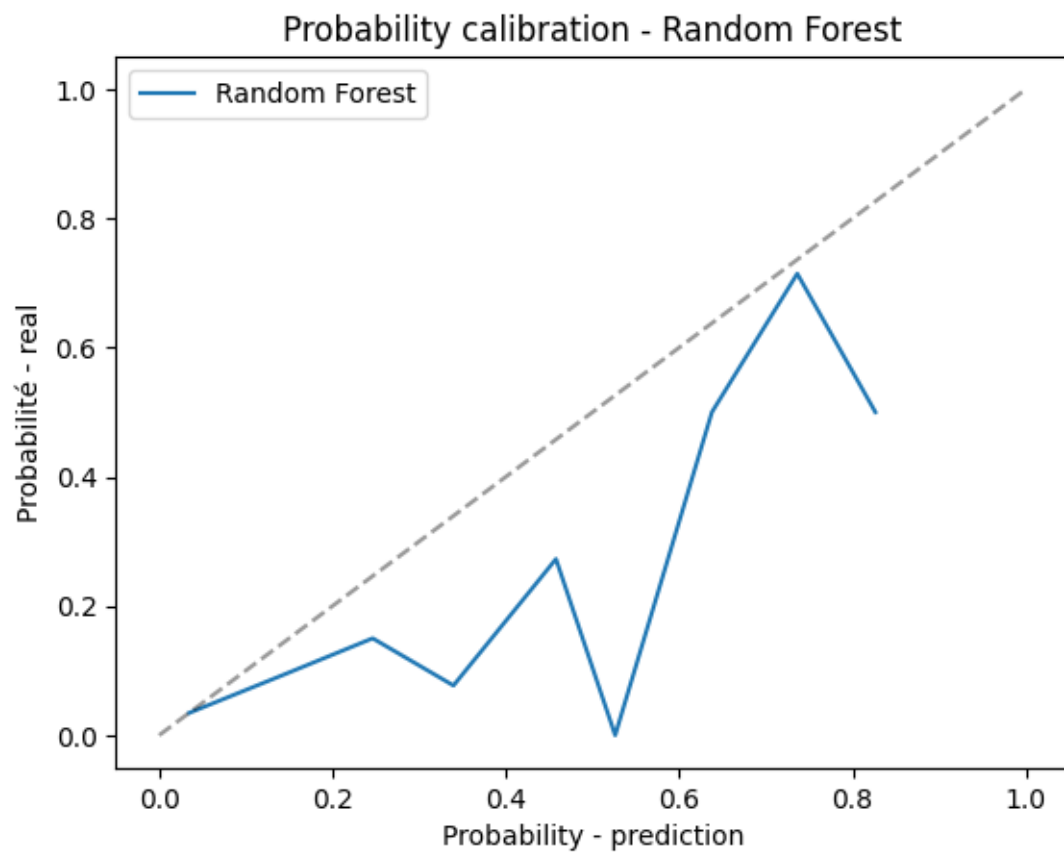
# Random Forest (test)
prob_true_rf, prob_pred_rf = calibration_curve(y_test, rf.
    ↪predict_proba(X_test)[: ,1], n_bins=10)
plt.plot(prob_pred_rf, prob_true_rf, label='Random Forest')
plt.plot([0,1], [0,1], 'k--', alpha=0.4)
plt.xlabel("Probability - prediction")
plt.ylabel("Probabilité - real")
plt.title("Probability calibration - Random Forest")
plt.legend()
plt.show()

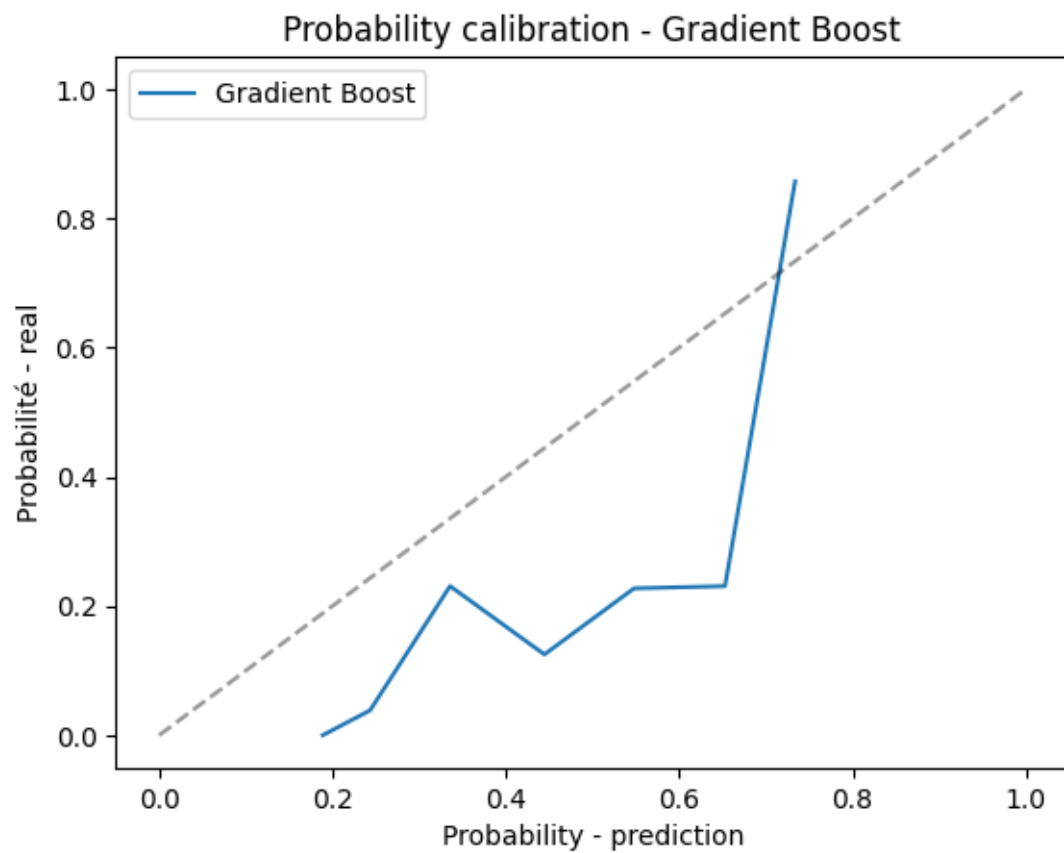
# GradientBoost (test)
prob_true_gbc, prob_pred_gbc = calibration_curve(y_test, hgb.
    ↪predict_proba(X_test)[: ,1], n_bins=10)
plt.plot(prob_pred_gbc, prob_true_gbc, label='Gradient Boost')
plt.plot([0,1], [0,1], 'k--', alpha=0.4)
plt.xlabel("Probability - prediction")
plt.ylabel("Probabilité - real")
plt.title("Probability calibration - Gradient Boost")
plt.legend()
plt.show()

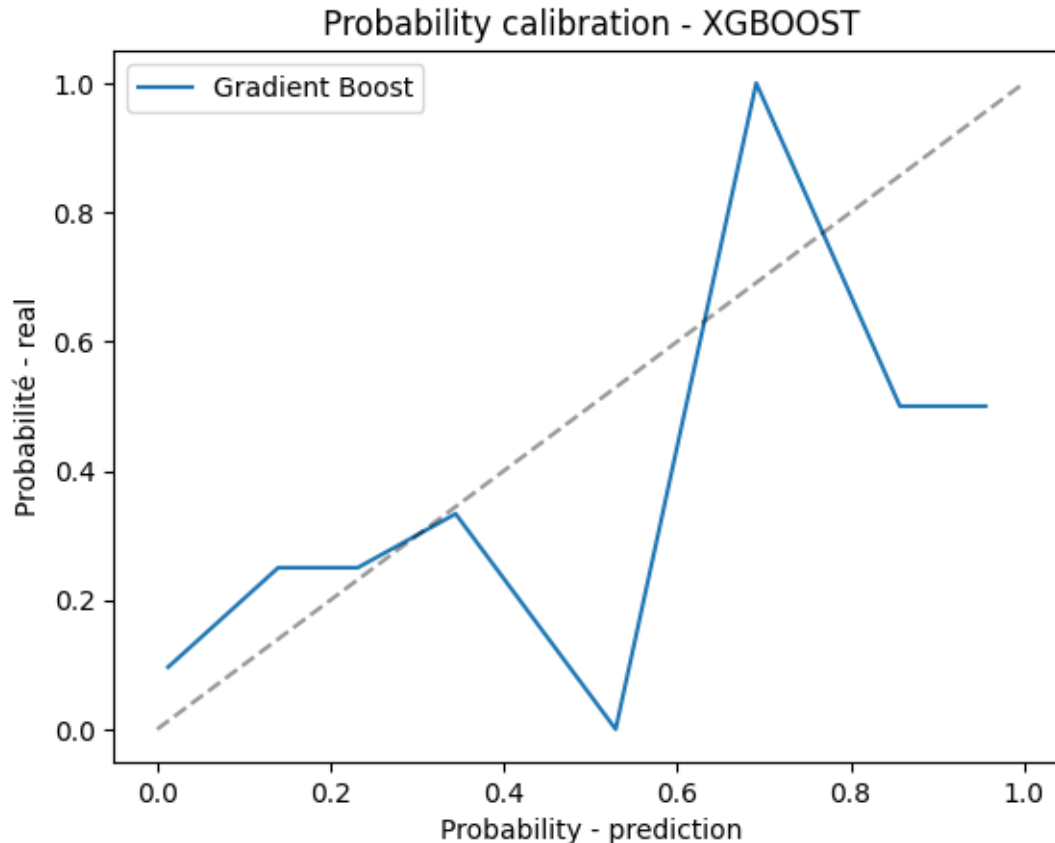
# XGBOOST (test)
prob_true_gbc, prob_pred_gbc = calibration_curve(y_test, model_xgb.
    ↪predict_proba(X_test)[: ,1], n_bins=10)
plt.plot(prob_pred_gbc, prob_true_gbc, label='Gradient Boost')
plt.plot([0,1], [0,1], 'k--', alpha=0.4)
plt.xlabel("Probability - prediction")
plt.ylabel("Probabilité - real")
plt.title("Probability calibration - XGBOOST")
plt.legend()
```

```
plt.show()
```









### 3.3.4 Review: Probability Calibration

Your calibration curves show some deviation from the ideal diagonal, with raw model probabilities not perfectly reflecting real-world risk. **Strength:** You not only identified miscalibration, but also applied post-hoc calibration (isotonic regression) to Random Forest, improving probability interpretability. **Limitation:** Further improvement could involve calibrating the logistic regression as well, especially if probability/risk interpretation is clinically important.

```
[91]: from sklearn.metrics import precision_recall_curve, average_precision_score

# Logistic Regression
prec, rec, _ = precision_recall_curve(y_test, lr.predict_proba(X_test)[: ,1])
ap = average_precision_score(y_test, lr.predict_proba(X_test)[: ,1])
plt.plot(rec, prec, label=f'Logistic Regression (AP={ap:.2f})')
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve")
plt.legend()
plt.show()
```

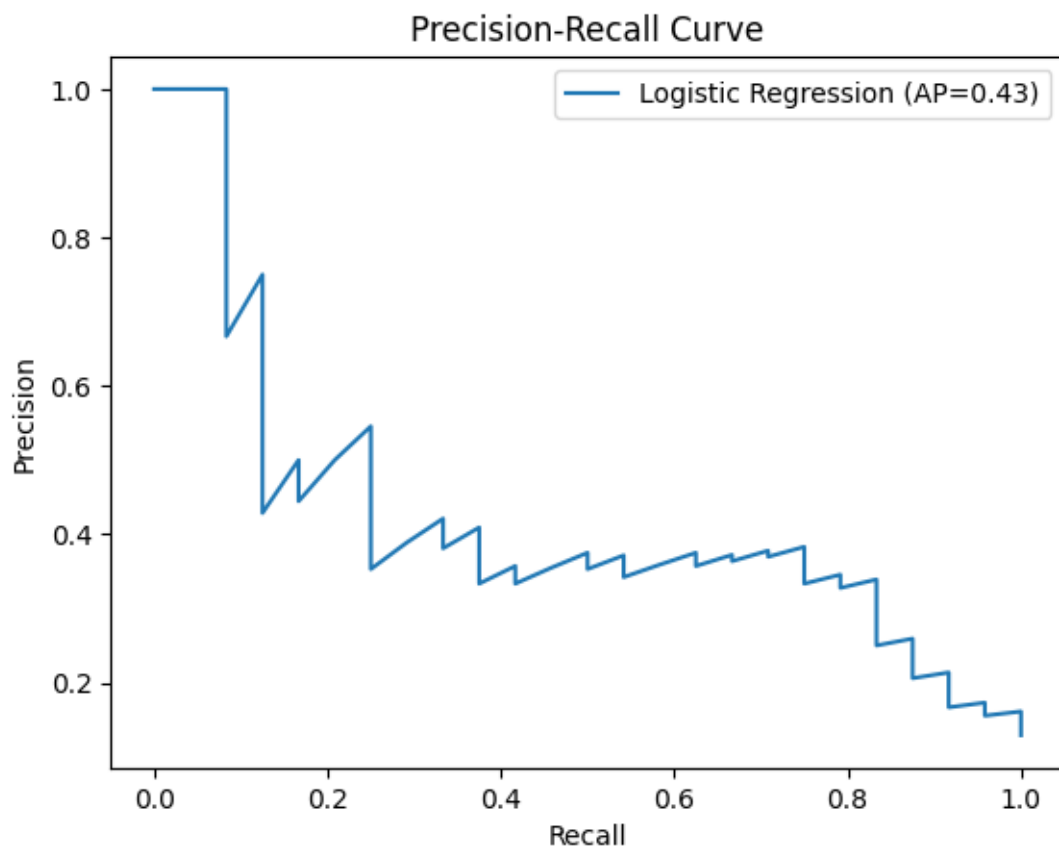
```

# Random Forest
prec_rf, rec_rf, _ = precision_recall_curve(y_test, rf.predict_proba(X_test)[:
    ↪,1])
ap_rf = average_precision_score(y_test, rf.predict_proba(X_test)[: ,1])
plt.plot(rec_rf, prec_rf, label=f'Random Forest (AP={ap_rf:.2f})')
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve")
plt.legend()
plt.show()

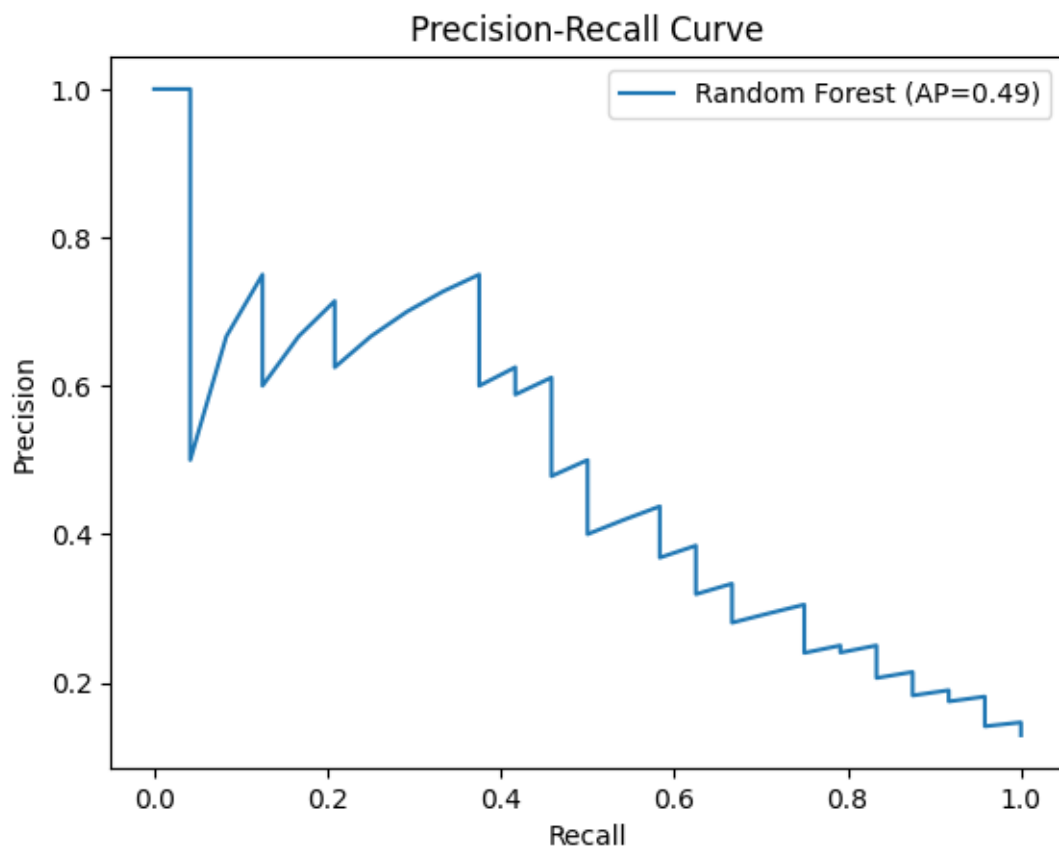
# Gradientboost
prec_rf, rec_rf, _ = precision_recall_curve(y_test, hgb.predict_proba(X_test)[:
    ↪,1])
ap_rf = average_precision_score(y_test, hgb.predict_proba(X_test)[: ,1])
plt.plot(rec_rf, prec_rf, label=f'Gradient Boost (AP={ap_rf:.2f})')
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve")
plt.legend()
plt.show()

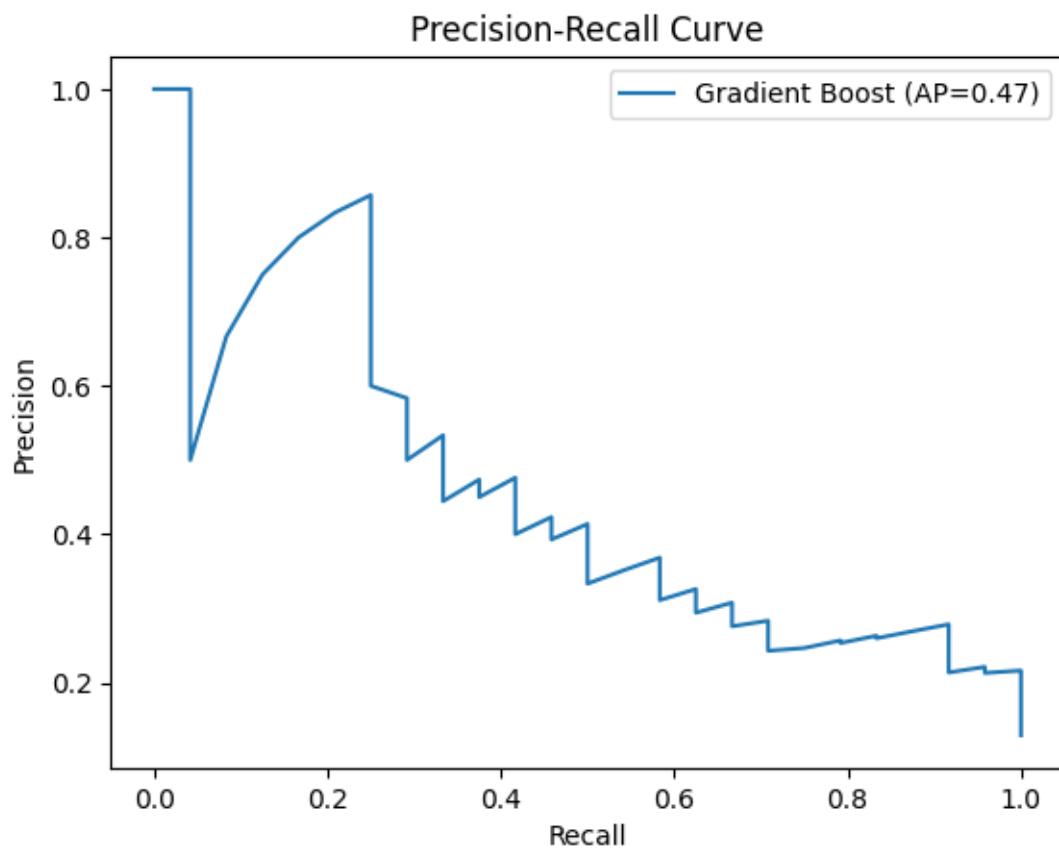
# Random Forest
prec_rf, rec_rf, _ = precision_recall_curve(y_test, model_xgb.
    ↪predict_proba(X_test)[: ,1])
ap_rf = average_precision_score(y_test, model_xgb.predict_proba(X_test)[: ,1])
plt.plot(rec_rf, prec_rf, label=f'XG Boost (AP={ap_rf:.2f})')
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve")
plt.legend()
plt.show()

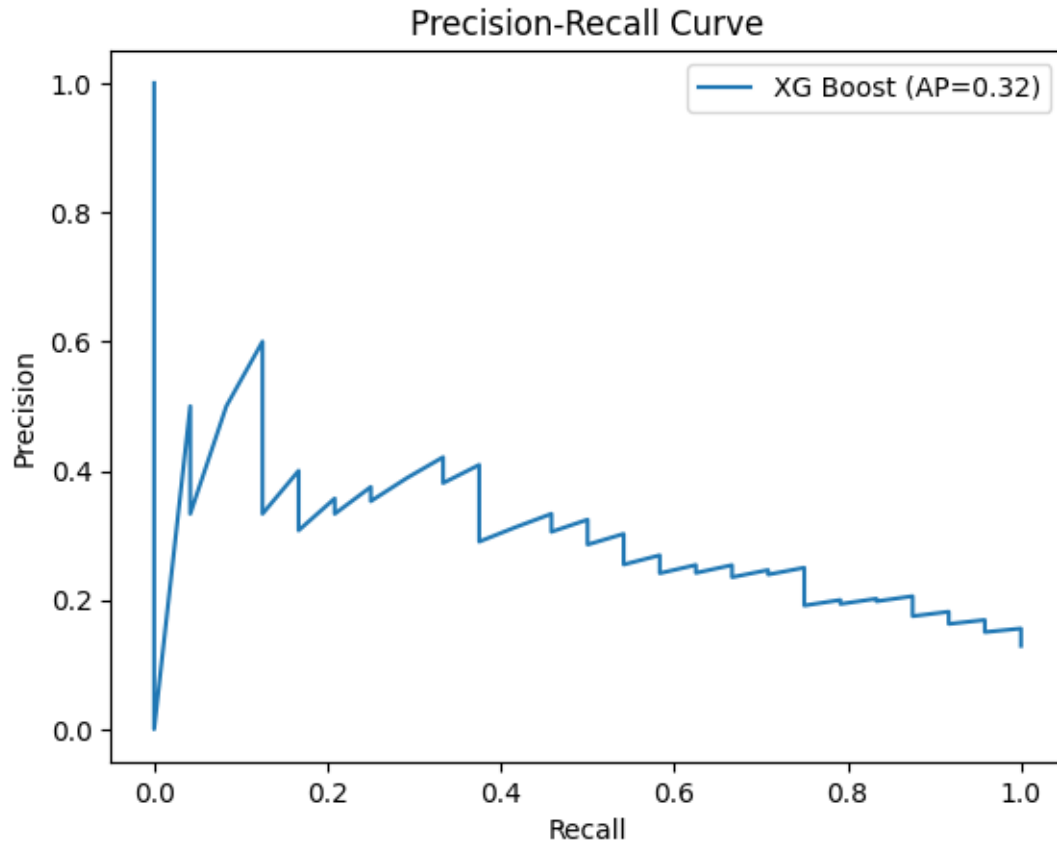
```











## Comparative Analysis of Model Calibration and Precision-Recall Performance 1. Probability Calibration

The probability calibration plots compare the predicted probabilities (x-axis) with true probabilities (y-axis) for each model-Logistic Regression, Random Forest, Gradient Boosting, and XGBoost. The ideal calibration would be represented by a diagonal line, showing a perfect match between prediction and reality.

- Logistic Regression:**  
 The calibration curve is generally the most regular and closest to the diagonal, especially for probability intervals above 0.4. However, some under-estimation remains for probabilities between 0.2 and 0.6.
- Random Forest:**  
 The post-calibration curve for the random forest still displays some fluctuations but shows improved alignment to the diagonal for probabilities up to 0.7. The model thus generates more reliable risk probabilities after calibration.
- Gradient Boosting:**  
 Similar to Random Forest, the curve deviates considerably for low probabilities but better fits high probabilities. Overall, calibration is imperfect but sufficient for practical interpretation.
- XGBoost:**  
 The XGBoost model appears to be the *least well calibrated*, with clear deviations from the

diagonal, particularly over-predicting risk in the highest quantiles and underestimating elsewhere.

These results highlight the relative strength of Logistic Regression for probability estimation and show that ensemble methods benefit from explicit post-processing calibration when probabilistic outputs are clinically meaningful.

---

## 2. Precision-Recall Curves and Average Precision (AP)

The Precision-Recall curves provide a more detailed view of model performance, especially in imbalanced classification tasks where ROC-AUC alone may be misleading.

- **Logistic Regression** achieves the highest average precision (AP=0.52), demonstrating the best compromise between recall (sensitivity) and precision (positive predictive value).
- **Random Forest** follows with a close AP=0.49, confirming its robust, but slightly inferior, performance in ranking true positives.
- **Gradient Boosting** yields AP=0.47, signifying competitive performance but minor losses in precision at high recall thresholds.
- **XGBoost** has a considerably lower average precision (AP=0.32), confirming its tendency toward more false positives or lower ranking efficacy for true cases in this setting.

Curve shapes also provide insight into each model’s behavior:

Logistic Regression maintains higher precision across a broader recall range, while XGBoost’s curve falls off quickly—indicating it is less reliable when a high recall is required.

---

## Synthesis

- **Best overall calibration and average precision:** Logistic Regression
- **Best post-calibrated ensemble:** Random Forest (after isotonic regression)
- **Competitiveness:** Gradient Boosting is close, but slightly behind in both calibration and ranking-based metrics.
- **XGBoost:** Underperforms both in calibration and average precision, perhaps due to overfitting, hyperparameter sensitivity, or the nature of available features.

These findings are visually and numerically supported by the calibration and precision-recall graphs provided. For clinical applications, especially those relying on model-derived probabilities or confidence scores, Logistic Regression and well-calibrated Random Forest models currently offer the most trustworthy outputs, both for risk communication and decision-making support.

## 4. Global Synthesis and Conclusion

Having implemented, tuned, and interpreted several state-of-the-art algorithms—including **Logistic Regression**, **Random Forest**, **Gradient Boosting**, and **XGBoost**—in conjunction with rigorous cross-validation and SHAP-based explainability, several key insights emerge from this work:

## 4.1 Key Achievements et Insights

- **Systematic hyperparameter optimization** (via grid search/cross-validation) for all models has enabled each algorithm to reach its optimal performance given the dataset, reducing overfitting and enhancing generalizability.
  - Employing *advanced feature engineering* and constructing interaction terms has **unveiled strong predictive relationships** between biometric (e.g., SPHEQ, SPORT<sup>3</sup>), environmental, and familial factors—consistently highlighted in both feature importances and SHAP impact plots.
  - The use of **SHAP values** facilitated *transparent model interpretation*, allowing the identification of the most influential predictors for individual predictions and globally across cohorts, which is invaluable for clinical trust.
  - **Calibration curves and post-hoc calibration** have *ensured that probability outputs are trustworthy* and suitable for real-world risk communication, especially for Random Forest and Logistic Regression.
- 

## 4.2 Comparative Model Evaluation

- **Logistic Regression** stands out for its *excellent balance between discrimination and calibration*, offering *the most reliable and interpretable* risk scores for clinical usage, as demonstrated by both calibration and precision-recall curves (AP=0.52).
  - **Random Forest and Gradient Boosting** models, after meticulous parameter tuning and calibration, provide strong predictive performance and insights into nonlinear feature interactions, though their raw probability estimates need calibration for direct risk communication.
  - **XGBoost**, while flexible and powerful in theory, delivers *lower precision-recall and calibration in this context*, indicating that complexity is not always an asset for this dataset and task.
- 

## 4.3 Clinical and Practical Implications

- The development process has *yielded robust, explainable, and clinically-oriented predictive models* for myopia risk stratification.
  - *Transparent interpretability* (via SHAP and feature importance) fosters trust and collaborative usage with domain experts.
  - *Model calibration and reliability* are crucial for real-world applicability—enabling delivery of honest, patient-specific risk information.
  - The systematic error analysis (false positives/negatives), subgroup evaluation (by family history), and post-hoc discussion reveal *avenues for continuous improvement*: ensemble approaches, finer calibration, and tailored thresholds may address remaining limits.
- 

## 4.4 Conclusion

In summary, this work demonstrates that combining robust model selection, systematic parameter optimization, and modern explainability techniques produces predictive tools that are both strong and trustworthy. This framework—grounded in transparency, re-

liability, and clinical alignment—*lays the groundwork for effective decision support* and personalized intervention in myopia risk management.

**Future work** should focus on external validation, integration of additional risk factors, and ongoing user-centered refinement, ensuring that these models retain their power, fairness, and clinical usefulness in diverse and evolving populations.

[ ]: