# CASE STUDIES

# -

# ANALYSIS OF A DIABETES HEALTH INDICATORS DATA SET

13/12/2023

82.05 - Análisis Predictivo - Final

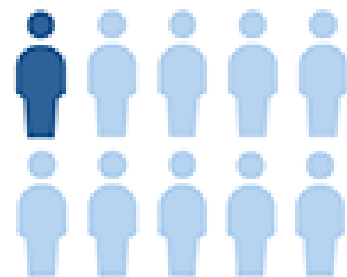Fanny LATRON (65998)

| | Diabetes_binary | HighBP | HighChol | CholCheck | BMI | Smoker | Stroke | HeartDiseaseorAttack | PhysActivity | Fruits | ... | AnyHealthcare | NoDocbcCost |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 1.0 | 0.0 | 1.0 | 26.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | ... | 1.0 | 0.0 |
| 1 | 0.0 | 1.0 | 1.0 | 1.0 | 26.0 | 1.0 | 1.0 | 0.0 | 0.0 | 1.0 | ... | 1.0 | 0.0 |
| 2 | 0.0 | 0.0 | 0.0 | 1.0 | 26.0 | 0.0 | 0.0 | 0.0 | 1.0 | 1.0 | ... | 1.0 | 0.0 |
| 3 | 0.0 | 1.0 | 1.0 | 1.0 | 28.0 | 1.0 | 0.0 | 0.0 | 1.0 | 1.0 | ... | 1.0 | 0.0 |
| 4 | 0.0 | 0.0 | 0.0 | 1.0 | 29.0 | 1.0 | 0.0 | 0.0 | 1.0 | 1.0 | ... | 1.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 70687 | 1.0 | 0.0 | 1.0 | 1.0 | 37.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 1.0 | 0.0 |
| 70688 | 1.0 | 0.0 | 1.0 | 1.0 | 29.0 | 1.0 | 0.0 | 1.0 | 0.0 | 1.0 | ... | 1.0 | 0.0 |
| 70689 | 1.0 | 1.0 | 1.0 | 1.0 | 25.0 | 0.0 | 0.0 | 1.0 | 0.0 | 1.0 | ... | 1.0 | 0.0 |
| 70690 | 1.0 | 1.0 | 1.0 | 1.0 | 18.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 1.0 | 0.0 |

**21 feature variables and 70692 survey responses**

# VARIABLES

**Variables :**

## Categorical :

- HighBP
- HighChol
- CholCheck
- Smoker
- Stroke
- HeartDiseaseorAttack
- PhysActivity
- Fruits
- Veggies

- HvyAlcoholConsump
- AnyHealthcare
- NoDocbcCost
- DiffWalk
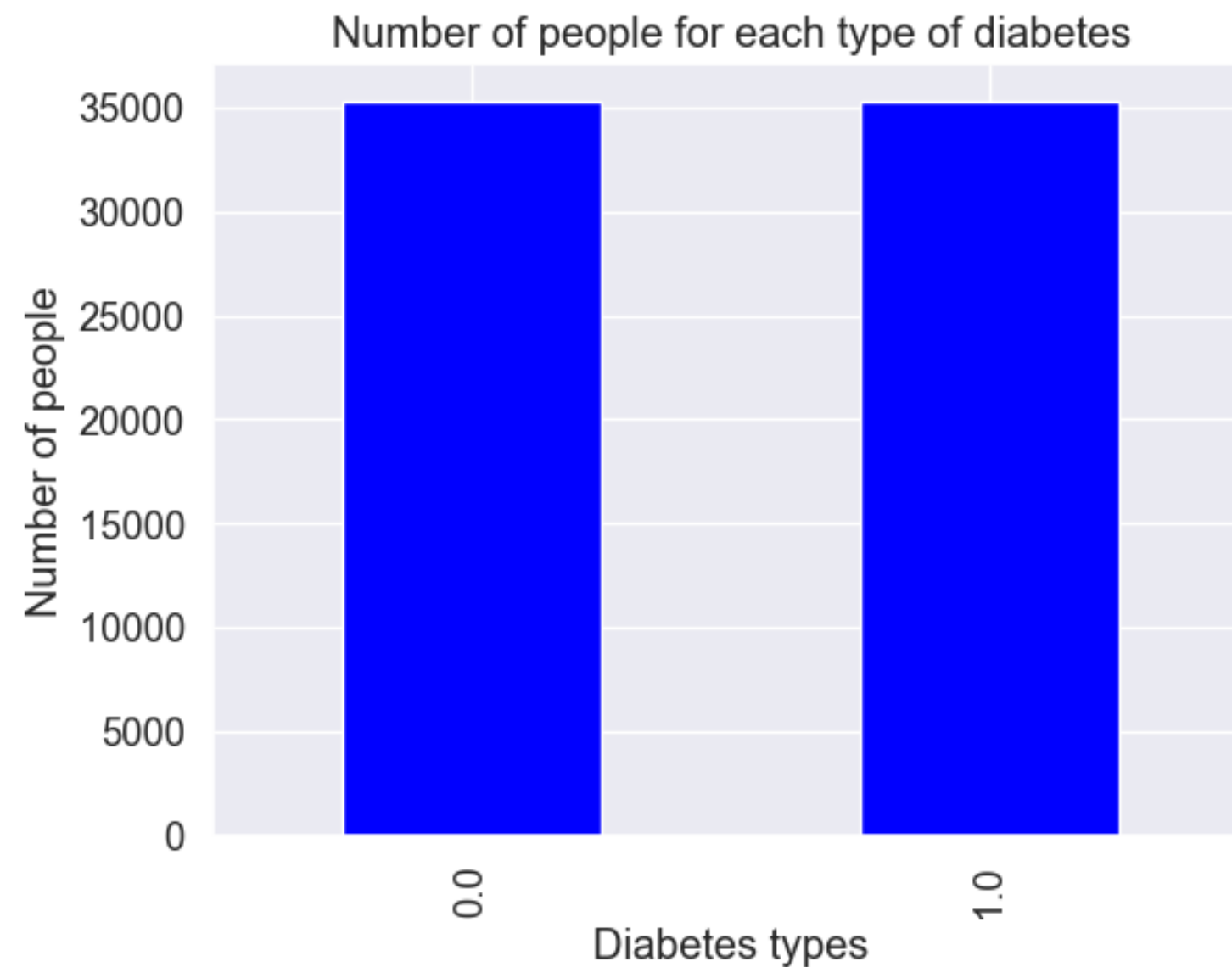- Sex
- GenHlth
- Age
- Education
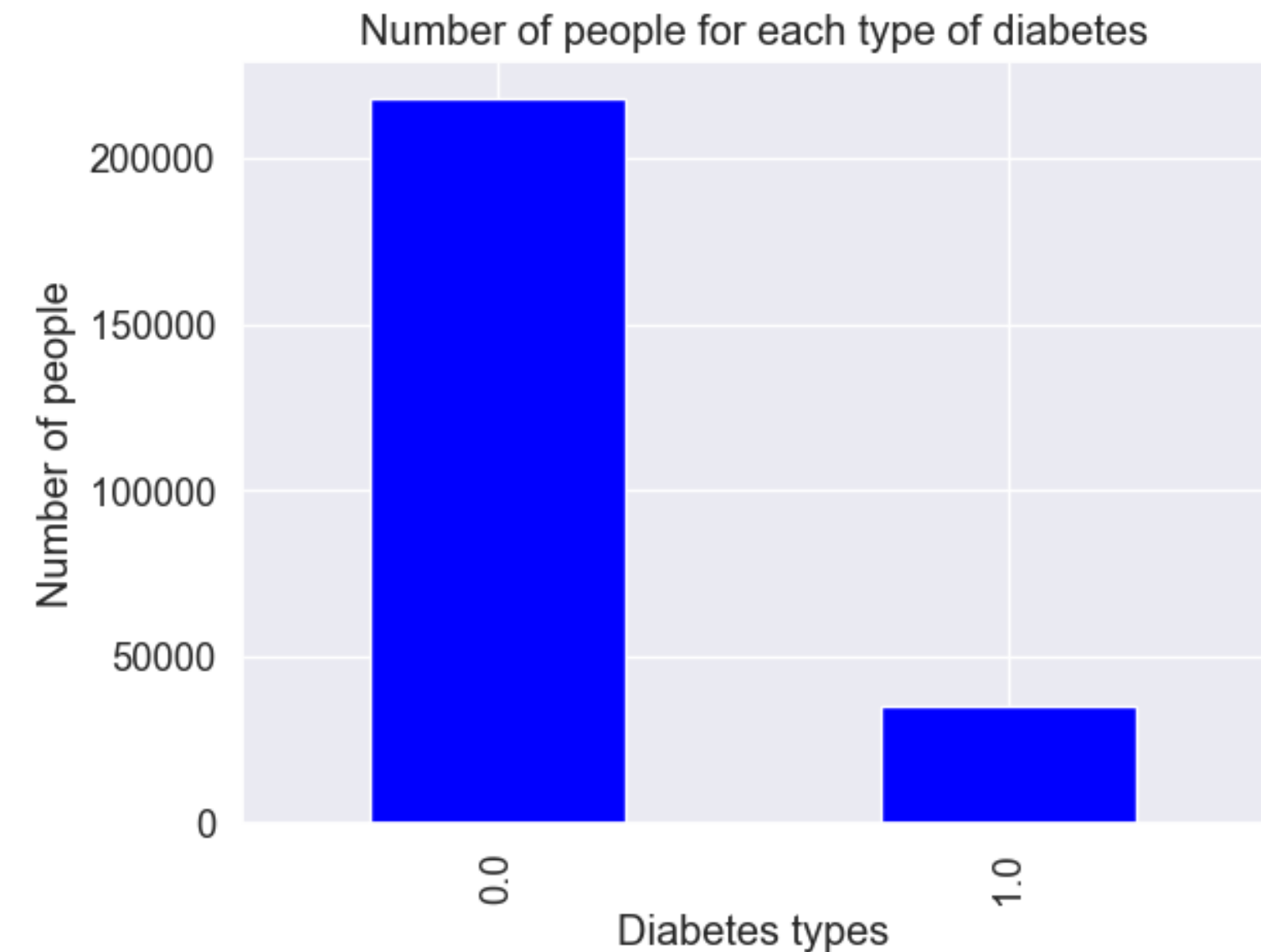- Income

## Numerical :

- MentHlth
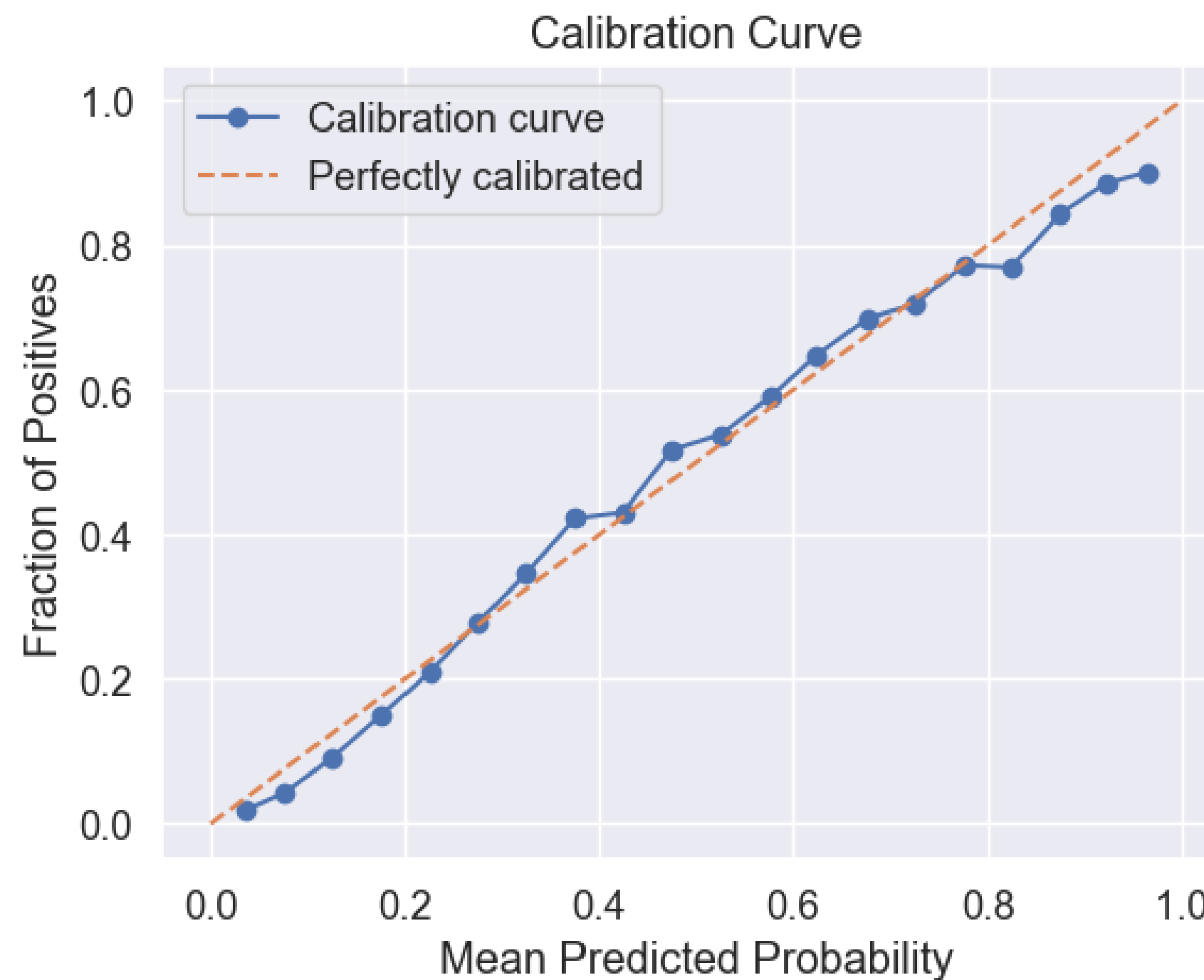- PhysHlth
- BMI

Cramers V correlation matrix

**● First Model : Logistic regression**

→ ROC-AUC Score: 0.8153
→ Brier Score: 0.1827



Calibration Curve

# Second Model : XGBoost
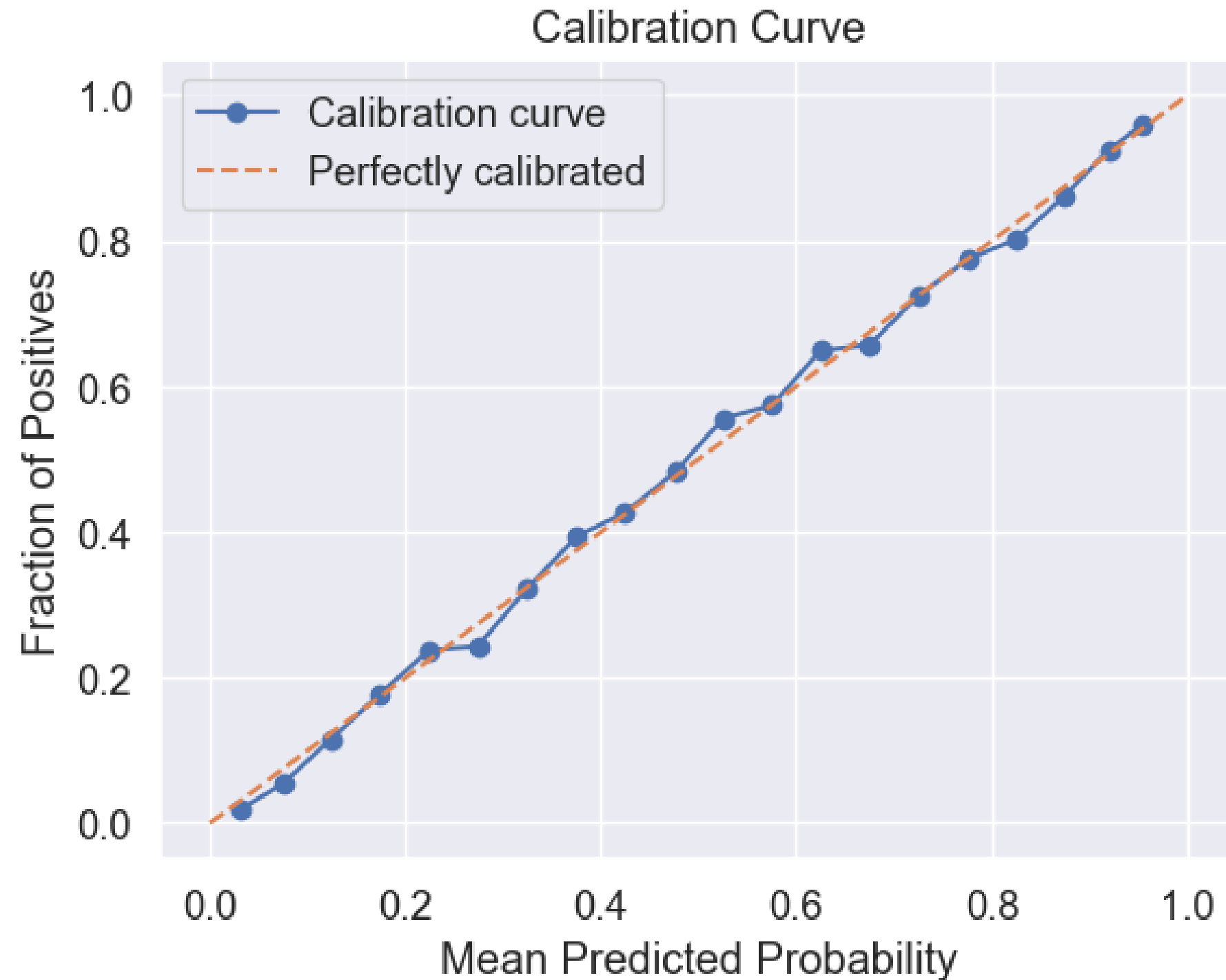
→ ROC-AUC Score: 0.8322
→ Brier Score: 0.1760
→ Parameters :
'learning_rate': 0.15,
'max_depth': 2,
'n_estimators': 200

# Other Model : SVM



Calibration Curve

# FINAL MODEL

**● Final Model : Random Forest**

→ ROC-AUC Score: 0.8461

→ Brier Score: 0.1739
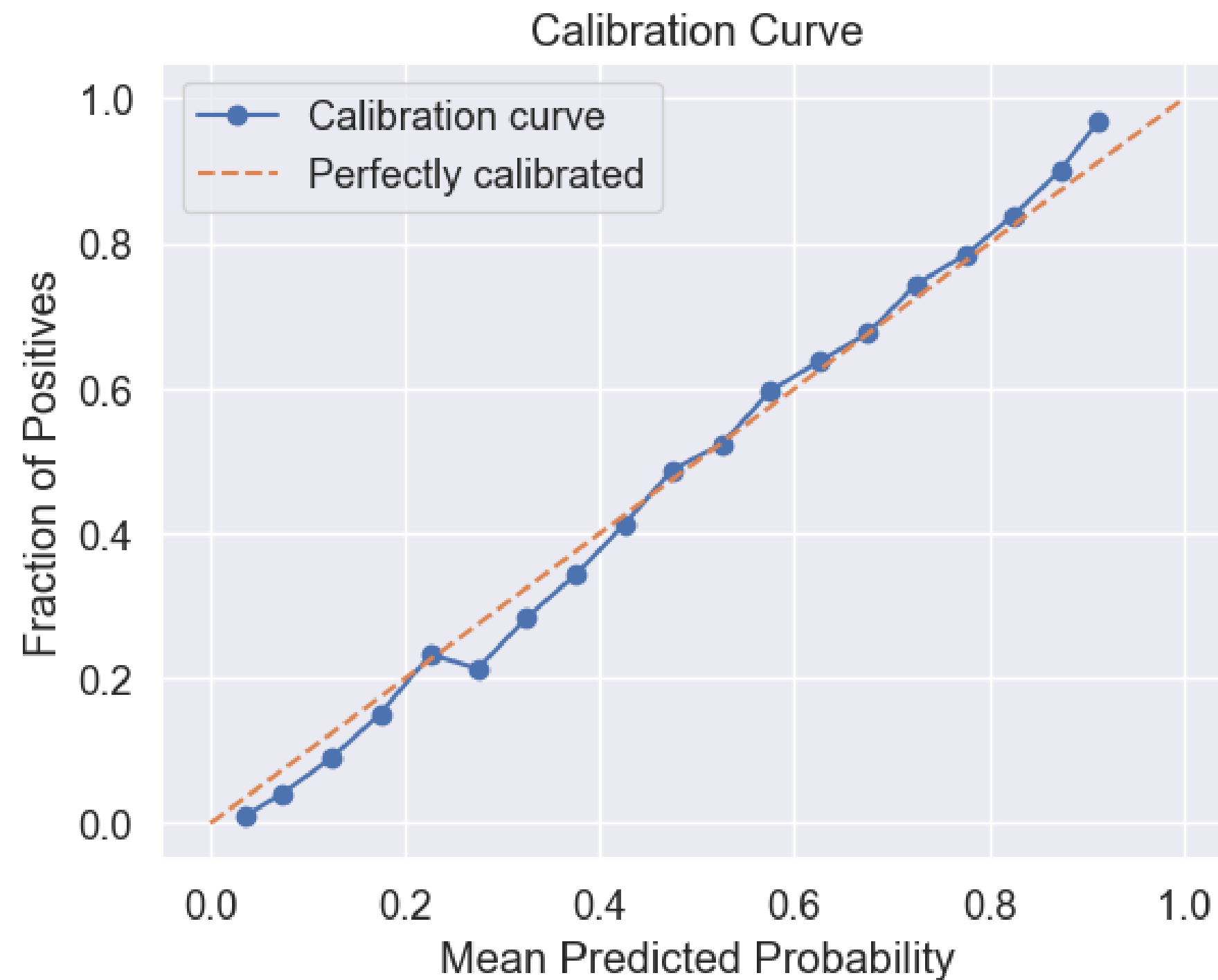
→ Parameters :
'max_depth': 12,
'min_samples_leaf': 3,
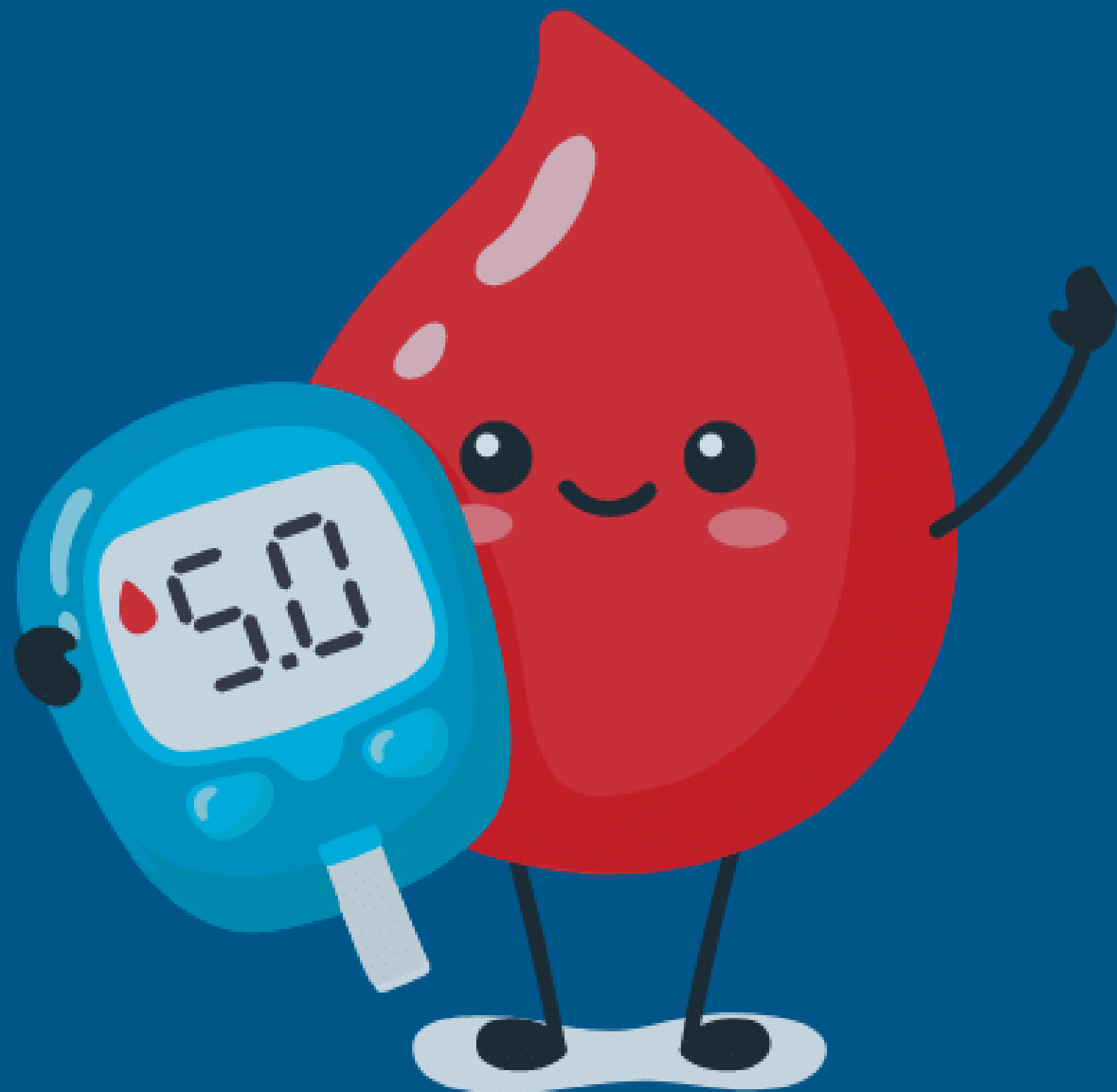'min_samples_split': 3,
'n_estimators': 210

Calibration Curve

# ITBA

# CONCLUSION

- **Could possibly find a better model but still efficient**

- **Low execution time**

- **Gives an indication close enough to the truth**

Thanks !