Statistical and Predictive analysis to Identify Risk Factors and Effects of Post COVID-19 Syndrome

Milad Leyli-abadi*, Jean-Patrick Brunet⁺ Sonja van Ockenburg⁺, Axel Tahmasebimoradi*

IRT SystemX, Palaiseau, France *
University Medical Center Groningen, Netherlands †







Introduction

Context and motivation

Data and statistical analysis

Predictive analysis of long COVID intensity

Conclusions and perspectives

Outline

Context & Motivation

In May 2023, the WHO declared the end of the COVID-19 global Public Health Emergency, marking a new endemic phase

However, COVID-19 continues to have lasting effects, notably Post-COVID-19 Condition (PCC) — persistent physical and cognitive symptoms following recovery

Recent global estimates suggest a **6–10% prevalence**, down from initial 10–20% WHO estimates

Identifying populations at risk of PCC is essential for early referral and optimized care pathways

Research Challenges & Data Opportunity

PCC characterization remains uncertain, with key challenges:

- Incomplete, retrospective data collected during the evolving crisis
- Focus on hospitalized patients, despite most PCC cases occurring in **non-hospitalized individuals**
- Limited information on **pre-existing conditions**, complicating symptom attribution

Lifelines Cohort offers a unique opportunity:

- 167,729 participants from Northern Netherlands across three generations
- COVID-19 branch (Apr 2020 Nov 2022):
 - 31 questionnaires, weekly to bi-monthly
 - 76,503 respondents, average 13.5 questionnaires each
- Rich longitudinal data enable analysis of pre-infection factors and PCC trajectories

Aim & Contributions

Research Question:

• Can pre-infection parameters predict the severity of Post-COVID-19 Condition?

Approach:

- Introduce Post-COVID-19 Symptom Intensity (PCSI) as a continuous measure of PCC severity
- Apply machine learning models to Lifelines data to predict PCSI and identify risk factors

Key Contributions:

- Statistical identification of influential factors associated with PCC
- Predictive modeling of PCSI using data-driven techniques
- Interpretation of variable impact for medical decision-making
- Development of an open-source Python package for reproducible ML analysis

Dataset Overview

Two main variable types:

- **Static variables:** Fixed individual attributes (e.g., age, sex, variant, income, smoking, chronic diseases, vaccination status, time between vaccination and infection)
- Dynamic variables: Symptom presence and intensity over time (before, during, and after infection)
 - Includes headache, cough, fever, breathing difficulties, muscle pain, loss of smell/taste, etc.

Data Challenges:

- Substantial **missing and inconsistent data** (common in questionnaire-based studies)
- Evolving questionnaire design during the epidemic → variable phrasing and coverage
- Required extensive standardization to build a uniform analytical dataset

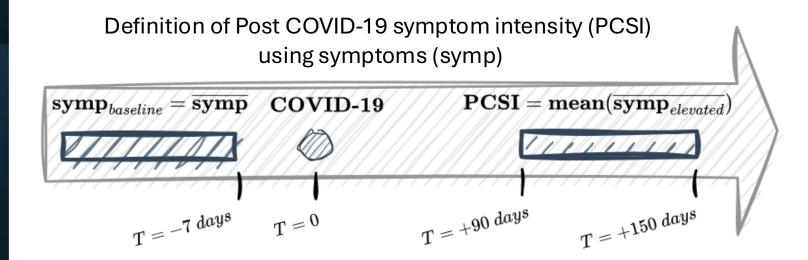
Defining Post-COVID-19 Symptom Intensity (PCSI)

Definition of Post-Covid Condition based on the WHO one

Symptoms lasting more than 3 months post infection lasting for at least 2 months with no other explanation.

Use of 10 symptoms checked to be related to PCC in a previous study

Extension to a continuous variable by deriving the Post Covid-19 symptom intensity score (PCSI)



Preprocessing Workflow & Sample Selection

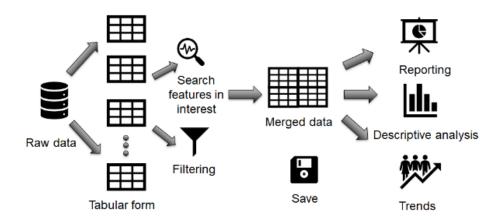
Data Integration:

- Raw questionnaire data organized by participant and date
- Merged across tables after cleaning and variable harmonization
- Feature extraction based on **pre-infection period** (steady-state hypothesis)

Final Study Sample:

- 4,657 participants included after filtering for sufficient data
- Comparable baseline characteristics between included and excluded groups
- Sex distribution:
 - Women = 73% of PCC cases vs. 64% in base dataset → higher PCC risk in women
 - Lower PCSI levels show a reduced female proportion

Data preparation strategy



Statistical Tests

Tests Used:

- Chi-square test → checks independence between categorical variables (e.g., vaccination & PCSI
- Cramer's V → measures strength of association (0 = weak, 1 = strong)

Vaccination & PCSI:

- Majority of **fully vaccinated** participants had **low PCSI** (88% scored 1–2)
- Chi-square: significant relationship (p < 0.05)
- Cramer's V = 0.072: weak association strength

Insight:

 Vaccination status is significantly associated with lower PCC intensity, but the strength of this association is relatively small in practical terms

 $(H_0:Independence\ between\ variables)$

 H_1 : Dependence between variables

VACCINIE	PC_INTENSITY					
VACCINE	1	2	3	4	5	Total
complete vaccin	2514	276	225	108	26	3149
	79.8 %	8.8 %	7.1 %	3.4 %	0.8 %	100 %
	71 %	54.9 %	59.7 %	54.5 %	70.3 %	67.6 %
no	1028	227	152	90	11	1508
	68.2 %	15.1 %	10.1 %	6 %	0.7 %	100 %
	29 %	45.1 %	40.3 %	45.5 %	29.7 %	32.4 %
Total	3542	503	377	198	37	4657
	76.1 %	10.8 %	8.1 %	4.3 %	0.8 %	100 %
	100 %	100 %	100 %	100 %	100 %	100 %

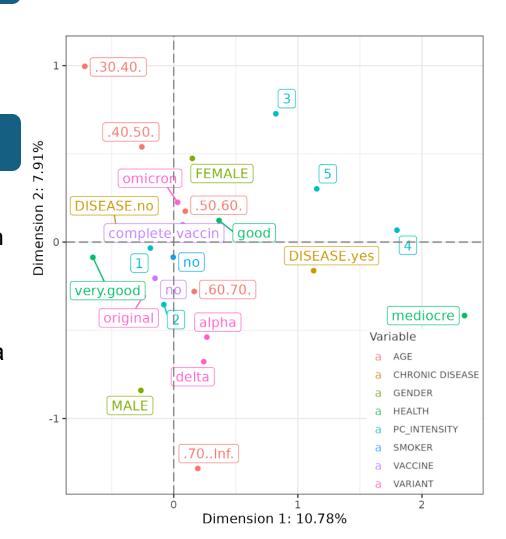
Multiple Correspondence Analysis (MCA)

Method

 MCA was applied to explore simultaneous relationships between multiple categorical variables

Key patterns

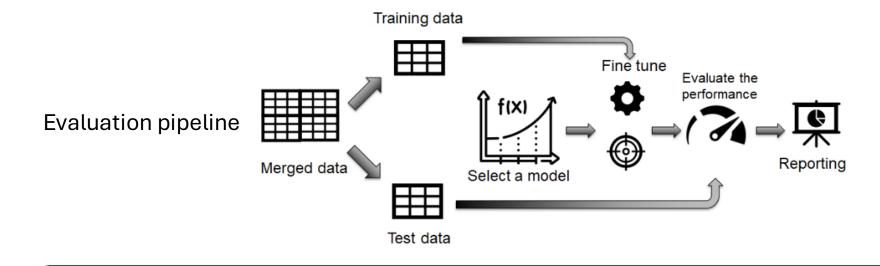
- High PCSI (5) clusters with chronic diseases, poor health, and female sex, indicating higher vulnerability in these groups
- SARS-CoV-2 variant shows weak association with PCC, suggesting a lower impact on long-term symptom intensity
- Participants with better general health are positioned closer to low PCSI levels, indicating a protective effect



Predictive analysis

The model $f: X \to y$

- With $X \in \mathbb{R}^p$ representing the explanatory variables (features) of dimension p
- And y representing the target variable corresponding to the long covid intensity with $y \in [1; 5]$



- To be robust to data variation during the training and evaluation , the cross-validation strategy is adopted
- The final results are reported using mean and standard deviation over multiple folds

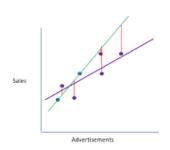
Cross validation

Training	Val	Test
60%	10%	/30%//
	/////	
	//	

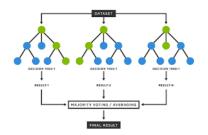
Methods and evaluation criteria

Evaluated methods

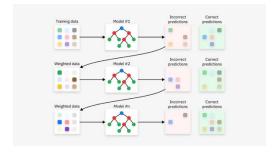
Linear Ridge Regression (LR)



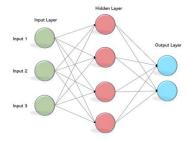
Random Forest (RF)



Gradient Boosting (GB)



Multi-layer Perceptron (MLP)



Evaluation criteria

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Mean Absolute Percentage Error (MAPE)
- Pearson Correlation

Results

- Best performance achieved with "All" features across all methods
- **Symptom-only features** yield similar results for LR, RF, GB; MLP benefits more from full feature sets
- MAE ≈ 0.60 for all models → acceptable error on a 1–5 PCSI scale
- MLP shows lowest MAPE (0.19), minimizing relative errors
- RF and GB yield the highest Pearson correlations, capturing complex linear interactions

Differences reflect **model** strengths:

- RF/GB: strong linear patterns → higher correlations
- MLP: non-linear modeling → better error minimization

		Evaluation criteria			
Methods	Features	MAE	MSE	MAPE	Pearson
LR	All	.61 \pm .01	.68 ± .02	.29 ± .01	(.56, 6e-70)
	Static	$.71 \pm .02$	$.91 \pm .05$	$.35 \pm .01$	(.28, 2e-16)
	Symptoms	$.62 \pm .02$	$.70\pm.04$	$.30 \pm .01$	(.57, 2e-69)
	Vaccination	$.81 \pm .02$	$.99\pm.05$	$.41 \pm .01$	NaN
RF	All	.60 \pm .01	.67 ± .02	.28 ± .01	(.58, 7e-73)
	Static	$.72 \pm .02$	$.93 \pm .05$	$.35 \pm .01$	(.26, 1e-15)
	Symptoms	.60 \pm .01	$.66 \pm .03$	$.28\pm.01$	(.57, 5e-72)
	Vaccination	$.79 \pm .02$	$.99 \pm .06$	$.39 \pm .01$	(.04, 1e-1))
GB	All	.61 \pm .01	.66 ± .01	.28 ± .01	(.57, 4e-74)
	Static	$.72 \pm .02$	$.90\pm.05$	$.35 \pm .01$	(.29, 7e-17)
	Symptoms	.61 \pm .01	$.68\pm.02$	$.28 \pm .01$	(.55, 8e-82)
	Vaccination	$.81 \pm .02$	$.99 \pm .06$	$.41 \pm .01$	(.05, 6e-1)
MLP	All	.45 ± .05	.90 ± .12	.19 ± .03	(.25, 3e-18)
	Static	$.87 \pm .18$	$1.4\pm.78$	$.43 \pm .07$	(.21, 4e-9)
	Symptoms	$1.76 \pm .11$	$.98\pm.38$	$.34 \pm .05$	(.43, 5e-33)
	Vaccination	$.80 \pm .03$	$1.03 \pm .05$.41 ± .03	(.04, 2e-1)

Interpretation

Linear Regression

Top 9 features identified via Linear Regression

- Coefficients show direction & strength of impact on PCSI
- **Positive predictors** (↑ PCC risk):
 - Loss of smell, headache, muscle pain
- Negative predictors (↓ PCC risk):
 - Fever, pain when breathing

Insights

• Acute symptoms differ in **predictive value** — some signal higher long-term risk, others appear protective

Coef	Variable	Coef
0.32	Pain when breathing	-0.58
0.28	Fever (38° or higher)	-0.27
0.27	Omicron variant	-0.26
0.23	Heaviness in arms/legs	-0.08
0.17	Very good health	-0.07
0.16	No chronic disease	-0.07
0.16	Age group	-0.06
0.16	Smoker	-0.05
0.14	Male	-0.03
	0.32 0.28 0.27 0.23 0.17 0.16 0.16	0.32 Pain when breathing 0.28 Fever (38° or higher) 0.27 Omicron variant 0.23 Heaviness in arms/legs 0.17 Very good health 0.16 No chronic disease 0.16 Age group 0.16 Smoker

Estimated coefficients of Linear Regression for prediction of post COVID-19 condition

Interpretation

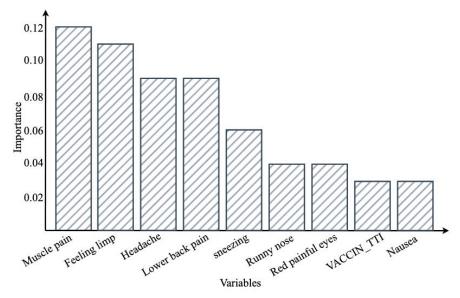
Random Forest

Top 10 features identified using **Random Forest**

Shows some overlap with Linear Ridge Regression results, but relative importance differs

Key Predictors:

- Muscle pain → most important predictor of PCSI
- Time between vaccination and infection (VACCIN_TTI) → significant impact
- Insight:
 - Vaccination timing may influence PCC risk and severity
 - RF highlights **non-linear interactions** not captured by linear models



Interpretation

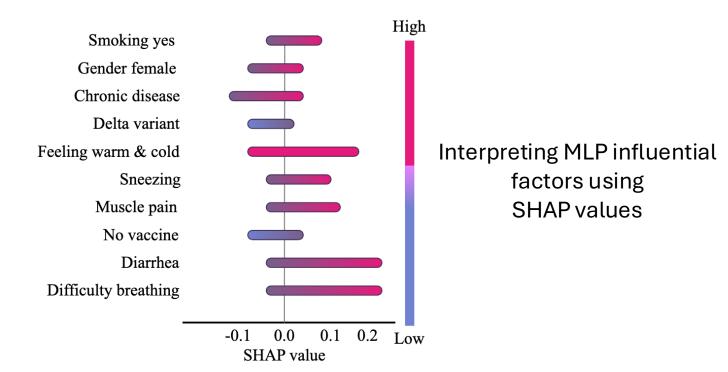
SHAP values

Top predictors for PCSI (via SHAP values):

- Positive impact (↑ PCSI): difficulty breathing, diarrhea, fluctuating body temperature, muscle pain, sneezing, smoking
- **Negative impact (PCSI):** prior vaccination, absence of chronic diseases

Other insights:

- Female sex associated with higher PCSI
- Highlights the complex interaction of symptoms, baseline health, and demographics in PCC risk



Conclusions and perspectives

Study Goal:

 Identify high-risk PCC profiles and predict Post-COVID Symptom Intensity (PCSI) using ML

Key Findings:

- **Higher risk:** women, patients with chronic diseases
- Important predictors: loss of smell, headache, muscle pain, vaccination timing
- Protective factors: absence of chronic diseases
- MLP slightly outperformed other models (lower MAPE)

Limitations:

- Steady-state assumption limits temporal analysis
- Dataset quality and completeness affect performance
- Models complement, **not replace**, clinical judgment

Societal Impact:

- PCC affects health, daily life, and productivity
- Predicting symptom intensity can guide early interventions and improve patient outcomes

References

[1] T. Lancet, The covid-19 pandemic in 2023: Far from over, 2023.

[2] A. V. Ballering, S. K. van Zon, T. C. olde Hartman, and J. G. Rosmalen, "Persistence of somatic symptoms after covid-19 in the netherlands: An observational cohort study," The Lancet, vol. 400, no. 10350, pp. 452–461, 2022

[3] F. Callard and E. Perego, "How and why patients made long covid," Social science & medicine, vol. 268, p. 113 426, 2021.

[4] https://www.who.int/europe/news-room/fact-sheets/item/post-COVID-19-condition.retrieved: September, 2025.

[5] C. E. Hastie et al., "Natural history of long-covid in a nationwide, population cohort study," Nature Communications, vol. 14, no. 1, p. 3504, 2023.

[6] R. Kessler, J. Philipp, J. Wilfer, and K. Kostev, "Predictive attributes for developing long covid—a study using machine learning and real-world data from primary care physicians in germany," Journal of Clinical Medicine, 2023.



Thank you for your attention!

Questions?