

Project Report

Marco Savastano
Carmine Vardaro

Information Engineering for Digital Medicine
Artificial Intelligence for Omics Data Analysis Course 2025-2026

Contents

1	Introduction	2
1.1	Background Clinico	2
1.2	La Metabolomica Untargeted	2
1.3	Problematiche Aperte	2
1.4	Scopo del Lavoro	2
2	Materials and Methods	2
2.1	Descrizione del Dataset	2
2.2	Quality Assessment & Data Cleaning	2
2.3	Strategie di Pre-processing (Iterativo)	2
2.4	Anomaly Detection	2
2.5	Feature Selection	2
2.6	Strategie di Data Fusion	2
2.7	Analisi Statistica e Machine Learning	2
2.8	Stack Tecnologico	2
3	Results and Discussion	3
3.1	Valutazione dell'Analisi Esplorativa (PCA)	3
3.2	Performance dei Modelli su Singoli Dataset (ESI+ / ESI-)	3
3.3	Risultati della Data Fusion	3
3.4	Interpretabilità e Biomarcatori (Feature Importance)	3
4	Conclusions	3
	References	3

Abstract

The short abstract (50-80 words) is intended to give the reader an overview of the work.

1 Introduction

1.1 Background Clinico

Breve panoramica sulla patologia (CHD) e importanza di trovare nuovi biomarcatori non invasivi.

1.2 La Metabolomica Untargeted

Perché la LC-MS è la scelta giusta qui (visione olistica del fenotipo).

1.3 Problematiche Aperte

Qui introduci il "problema" del tuo progetto: la complessità dei dati, la necessità di integrare due modi di ionizzazione (ESI+/ESI-) e la scelta del miglior metodo di preprocessing (non esiste una "ricetta unica").

1.4 Scopo del Lavoro

Valutare e ottimizzare un workflow chemiometrico (dal preprocessing alla Data Fusion) per distinguere soggetti Sani vs Patologici e identificare le feature biologicamente rilevanti.

2 Materials and Methods

2.1 Descrizione del Dataset

Origine dei dati (LC-MS). I due blocchi: ESI+ e ESI-. Composizione delle classi (Controlli vs CHD).

2.2 Quality Assessment & Data Cleaning

Descrizione dei QC e dei Replicati Tecnici. **Nota:** Spiegare chiaramente che QC e duplicati sono stati utilizzati per una valutazione iniziale della stabilità ma rimossi dal dataset finale di modeling (es. numero esiguo, impossibilità di mediare).

2.3 Strategie di Pre-processing (Iterativo)

Elenco delle tecniche testate: Normalizzazione, Trasformazione Logaritmica, Scaling (Autoscaling).

2.4 Anomaly Detection

Descrizione del metodo usato per identificare e rimuovere gli outlier (fondamentale per la pulizia del dato).

2.5 Feature Selection

Metodi statistici/algoritmici utilizzati per ridurre la dimensionalità e rimuovere il rumore prima del modeling.

2.6 Strategie di Data Fusion

Definizione degli approcci: Low-Level (concatenazione semplice) vs High-Level (o altri approcci). L'obiettivo è sfruttare la complementarità ESI+/ESI-.

2.7 Analisi Statistica e Machine Learning

- **Unsupervised:** PCA (per l'esplorazione).
- **Supervised:** PLS-DA, SVM, Random Forest, Logistic Regression.
- **Validazione:** Descrizione rigorosa dello split Training Set vs Validation Set (o Cross-Validation) per evitare l'overfitting.

2.8 Stack Tecnologico

Breve paragrafo sulle librerie Python utilizzate (Pandas, Scikit-learn, ecc.) per garantire la riproducibilità.

3 Results and Discussion

3.1 Valutazione dell'Analisi Esplorativa (PCA)

Visualizzazione dei dati ESI+ e ESI- separati. Valutazione degli Outlier (prima e dopo la rimozione). Confronto delle tecniche di scaling (es. efficacia dell'autoscaling).

3.2 Performance dei Modelli su Singoli Dataset (ESI+ / ESI-)

Confronto delle metriche (Accuratezza, Specificità, Sensibilità) tra PLS-DA, SVM, RF, LR. Quale modello performa meglio sui dati positivi? E sui negativi?

3.3 Risultati della Data Fusion

La fusione dei dati ha migliorato la classificazione rispetto ai dataset singoli?

3.4 Interpretabilità e Biomarcatori (Feature Importance)

Analisi delle Feature Importances e Analisi Univariata (Volcano Plot). Identificazione/interpretazione biologica dei top-metaboliti.

4 Conclusions

Sintesi del miglior workflow identificato. Considerazioni sull'interpretabilità biologica e limiti dello studio (es. numero di campioni, assenza di validazione esterna).

References

- [1] J. Hagenauer, E. Offer, and L. Papke. Iterative decoding of binary block and convolutional codes. *IEEE Trans. Inform. Theory*, vol. 42, no. 2, pp. 429—445, Mar. 1996.
- [2] T. Mayer, H. Jenkac, and J. Hagenauer. Turbo base-station cooperation for intercell interference cancellation. *IEEE Int. Conf. Commun. (ICC)*, Istanbul, Turkey, pp. 356–361, June 2006.
- [3] J. G. Proakis. *Digital Communications*. McGraw-Hill Book Co., New York, USA, 3rd edition, 1995.
- [4] F. R. Kschischang. Giving a talk: Guidelines for the Preparation and Presentation of Technical Seminars. <http://www.comm.toronto.edu/frank/guide/guide.pdf>.
- [5] IEEE Transactions \LaTeX and Microsoft Word Style Files. <http://www.ieee.org/web/publications/authors/transjnl/index.html>

TABLE I: Simulation Parameters

Information message length	$k = 16000$ bit
Radio segment size	$b = 160$ bit
Rate of component codes	$R_{cc} = 1/3$
Polynomial of component encoders	$[1, 33/37, 25/37]_8$

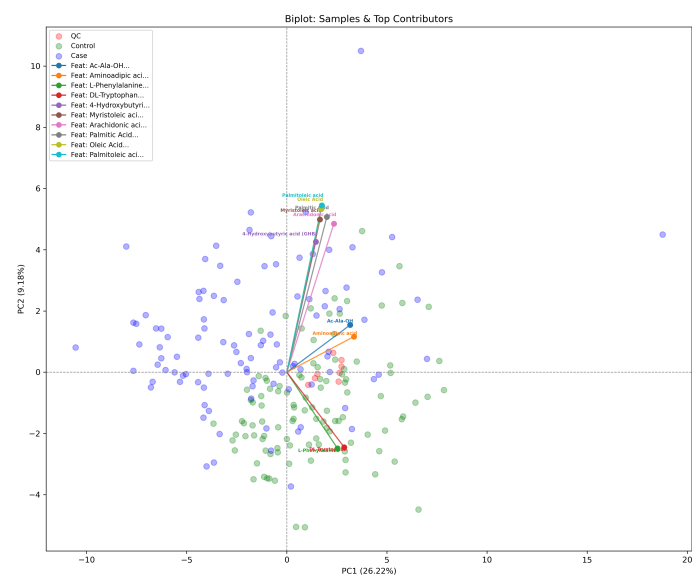


Fig. 1: Simulation results on the AWGN channel. Average throughput k/n vs E_s/N_0 .