



UNIVERSITÀ
CUSANO

UNIVERSITÀ TELEMATICA UNICUSANO

ACTIVITY 3

Machine Learning

Simone Arcari, IN32000132
22 november 2001

Capitolo 0: Contents

1	Introduzione	2
1.1	Obiettivi	2
2	Architettura del Sistema	3
2.1	Diagramma dell'architettura	4
3	Implementazione	5
3.1	requirements.txt	5
3.2	run.sh	5
3.3	src/core	6
3.3.1	Etivity3.py	6
3.3.2	Tee.py	6
3.4	src/gui	6
3.4.1	Etivity3Window.py	7
3.5	src/main.py	7
4	Interfaccia Utente	8
5	Esecuzione del Sistema	9
5.1	Flusso Operativo	9
5.1.1	Selezione del Dataset	10
5.1.2	Configurazione dell'Analisi	10
5.1.3	Esecuzione dell'Analisi	10
5.1.4	Visualizzazione dei Risultati	10
5.2	Interazione Avanzata	10
5.3	Output e Risultati	11

Capitolo 1: Introduzione

Il progetto è una soluzione software sviluppata in Python che permette di eseguire un'analisi automatica e completa di dataset. Lo strumento è progettato per semplificare il processo di analisi dei dati, offrendo funzionalità che spaziano dall'esplorazione iniziale dei dati fino all'applicazione di algoritmi di machine learning supervisionato e non supervisionato. I risultati vengono mostrati all'utente tramite una GUI intuitiva realizzata con PyQt5 e appositi grafici per una migliore interpretazione dei dati.

1.1 Obiettivi

- **Automatizzare** l'analisi dei dati riducendo il tempo necessario per l'elaborazione manuale
- **Classificare automaticamente** le variabili in base alla loro natura:
 - Numeriche continue/discrete
 - Categorie nominali/ordinali
- **Preprocessare** i dati in modo ottimale per l'analisi successiva
- **Eseguire analisi esplorative** per comprendere:
 - Distribuzione dei dati
 - Relazioni tra le variabili
- **Applicare tecniche di clustering** per l'analisi non supervisionata:
 - K-Means
 - Clustering gerarchico
- **Implementare modelli** per l'analisi supervisionata:
 - Classificazione (Regressione logistica, Albero decisionale)
 - Regressione (Regressione Lineare, Albero Regressore)
- **Generare visualizzazioni** chiare e interpretabili per facilitare la comprensione dei risultati

Capitolo 2: Architettura del Sistema

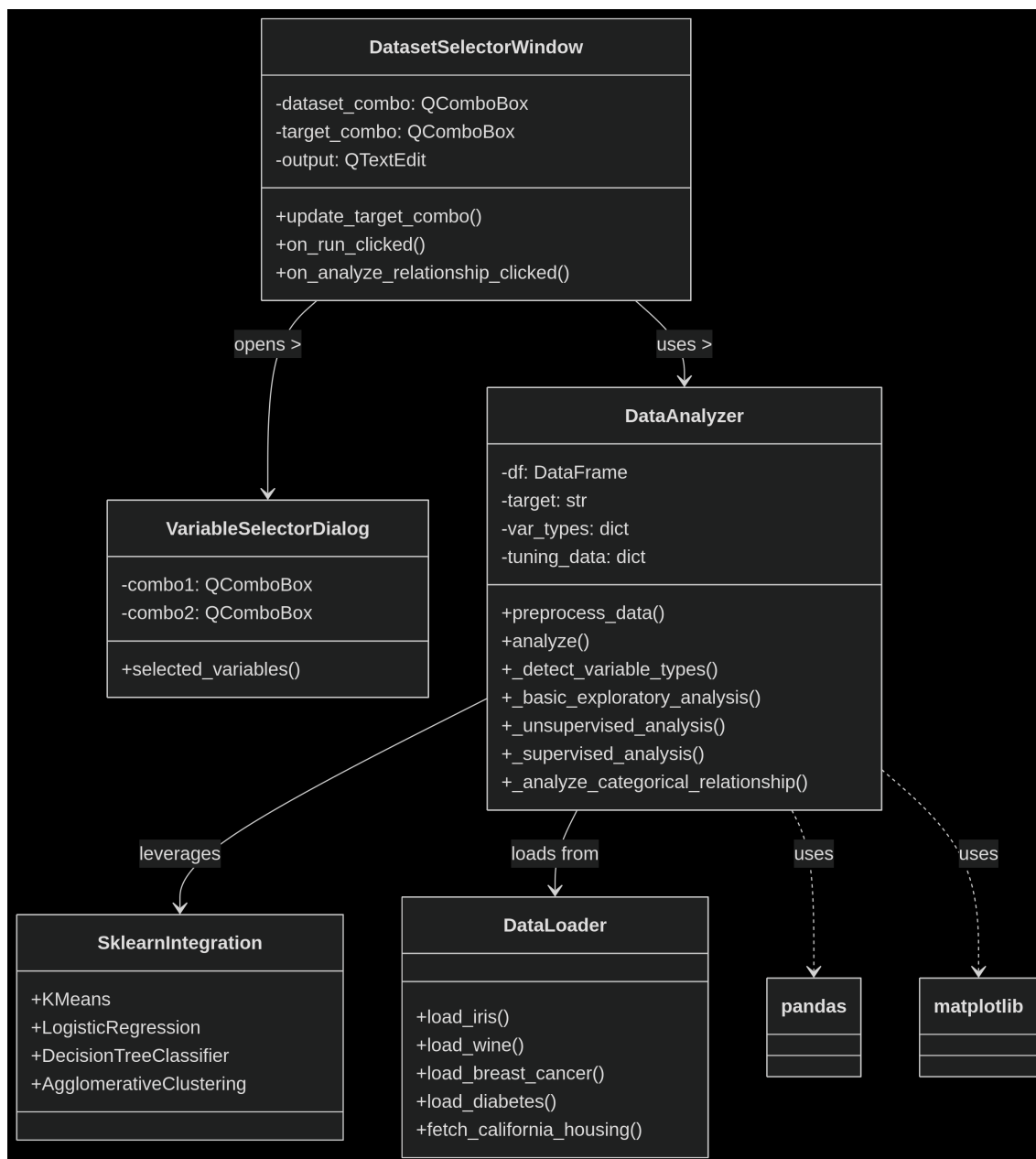
L'architettura del sistema si articola in tre livelli principali che dialogano per trasformare dati grezzi in risultati analitici. Il frontend, sviluppato con PyQt5, costituisce l'interfaccia utente attraverso cui gli utenti interagiscono con il sistema. La finestra principale `DatasetSelectorWindow` funge da hub centrale, permettendo la selezione dei dataset disponibili e l'impostazione dei parametri di analisi. Quando l'utente richiede un'analisi tra due variabili (Test Chi-Quadro), il sistema attiva la finestra `VariableSelectorDialog`, specializzata nella scelta delle coppie di variabili da esaminare.

Al cuore del sistema risiede il componente **DataAnalyzer**, il vero motore analitico che orchestra l'intero flusso di elaborazione. Questo modulo inizia classificando automaticamente le variabili rilevate, distinguendo tra tipologie numeriche continue, discrete e categoriche. Attraverso un processo di preprocessing, standardizza e codifica i dati, preparandoli per le successive fasi di analisi. La vera potenza del progetto si manifesta nella capacità di selezionare automaticamente la metodologia analitica più appropriata in base alla presenza o assenza di una variabile target, attivando rispettivamente analisi supervisionate o non supervisionate.

L'integrazione con l'ecosistema scikit-learn avviene attraverso un layer specializzato che include i principali algoritmi di machine learning. Per il clustering non supervisionato, il sistema può avvalersi sia dell'approccio K-Means che di quello gerarchico, mentre per le analisi supervisionate dispone sia di modelli di classificazione che di regressione. Un modulo dedicato gestisce il caricamento dei dataset predefiniti per l'analisi.

Il sistema completa la sua architettura con un meccanismo di visualizzazione che trasforma i risultati analitici in rappresentazioni grafiche, dalle heatmap di correlazione ai dendrogrammi, fino ai grafici a dispersione per la visualizzazione dei cluster.

2.1 Diagramma dell'architettura



Capitolo 3: Implementazione

Il progetto è organizzato in due macro-parti:

- **Core** (src/core): tutto il codice di calcolo e logica statistica, indipendente da GUI.
- **GUI** (src/gui + main.py): interfaccia utente PyQt5 per eseguire e visualizzare i risultati in modo interattivo.

Il progetto si sviluppa secondo la seguente alberatura:

```
ml-etivity3
├── requirements.txt
├── run.sh
├── src
│   ├── core
│   │   ├── data
│   │   │   └── tuning.py
│   │   ├── Etivity3.py
│   │   └── Tee.py
│   ├── gui
│   │   ├── Etivity3Window.py
│   └── main.py
```

3.1 requirements.txt

Elenca tutte le dipendenze Python (PyQt5, pandas, numpy, matplotlib, seaborn, ecc.) che verranno installate nell'ambiente virtuale.

3.2 run.sh

Script bash per avviare il programma, eseguendo 3 step principali:

- Crea (se necessario) e attiva un ambiente virtuale Python
- Installa le dipendenze da requirements.txt nell'ambiente virtuale
- Avvia main.py

3.3 src/core

Qui risiede tutta la logica di calcolo, indipendente dall'interfaccia grafica:

3.3.1 Etivity3.py

Implementa la classe `DataAnalyzer`, il cuore analitico del sistema che fornisce:

- **Inizializzazione:**
 - Caricamento dataset da file CSV o DataFrame esistente
 - Rilevamento automatico del tipo delle variabili (numeriche/categoriche)
 - Configurazione parametri di tuning per gli algoritmi
- **Preprocessing:**
 - Codifica variabili categoriche con `LabelEncoder`
 - Normalizzazione dati numerici con `StandardScaler`
- **Analisi Automatica:**
 - Statistiche descrittive e matrici di correlazione
 - Analisi non supervisionata (clustering K-Means e gerarchico)
 - Analisi supervisionata (regressione logistica, alberi decisionali, regressione lineare e albero regressore)
 - Test Chi-Quadro per relazioni tra variabili categoriche
- **Visualizzazione:**
 - Generazione automatica di grafici esplicativi
 - Heatmap per correlazioni
 - Dendrogrammi per clustering gerarchico

La classe è progettata per essere sia autonoma che integrabile con l'interfaccia grafica, esponendo tutti i risultati attraverso attributi di classe e metodi pubblici.

3.3.2 Tee.py

Implementa la classe `Tee`, che duplica ogni chiamata a `write()` su più stream. Usata assieme a `contextlib.redirect_stdout` per inviare contemporaneamente l'output a terminale e ad uno stream di stringhe.

3.4 src/gui

Qui risiede l'interfaccia grafica realizzata con PyQt5:

3.4.1 Etivity3Window.py

Implementa la finestra grafica con cui l'utente interagisce per eseguire l'analisi del dataset selezionato.

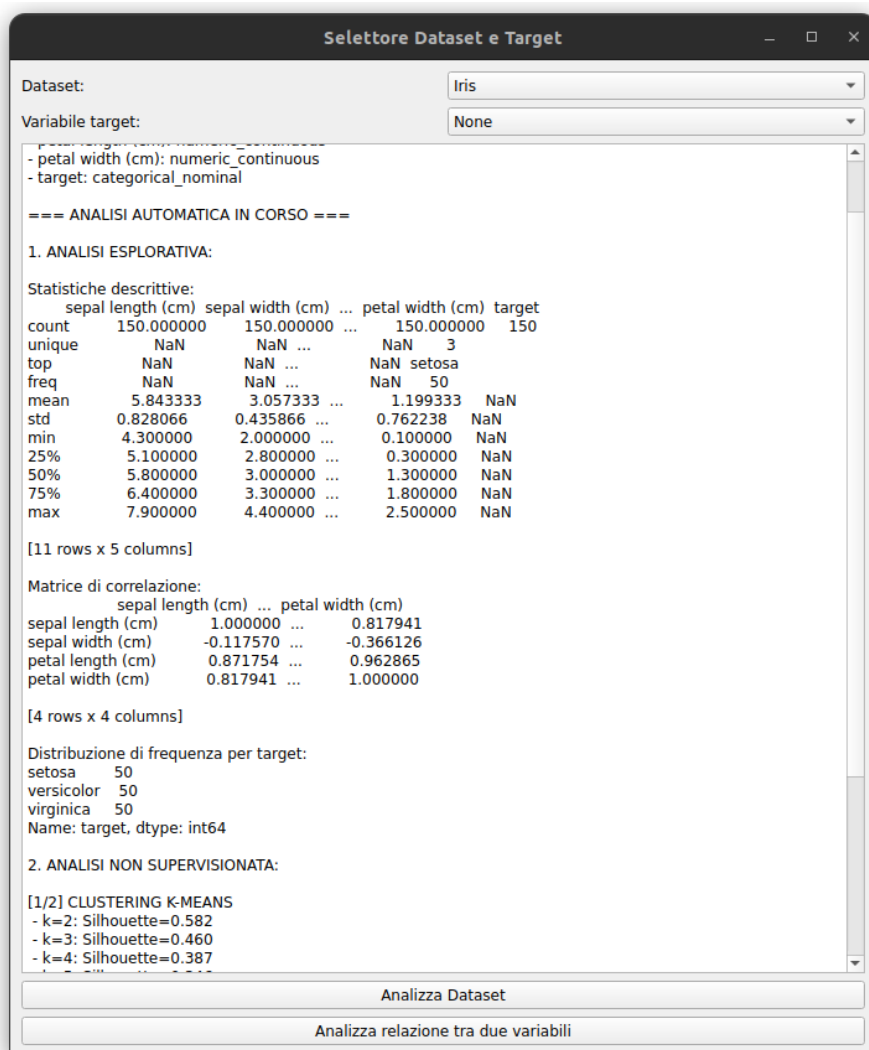
3.5 src/main.py

È il punto di ingresso dell'applicazione, chiamato dallo script di avvio `run.sh`. Il suo scopo è quello di avviare l'event loop per utilizzare l'interfaccia grafica Qt e gestire il segnale Ctrl+C per permettere all'utente di chiudere il programma da terminale.

Capitolo 4: Interfaccia Utente

L'interfaccia utente è stata sviluppata utilizzando PyQt5. La finestra principale permette all'utente di selezionare uno dei dataset preimpostati tramite un menu a tendina. Premendo il tasto *Analizza Dataset* inizia l'elaborazione sul dataset selezionato. I risultati vengono mostrati in un'area di testo non modificabile e tramite grafici a schermo.

Nota: se l'utente seleziona una variabile target tramite il secondo menu a tendina, verrà eseguita un'analisi supervisionata; in caso contrario, si esegue un'analisi non supervisionata.



Capitolo 5: Esecuzione del Sistema

L'applicazione viene avviata eseguendo lo script principale `run.sh`, che attiva l'interfaccia grafica sviluppata con PyQt5.

```

simone@simone:~/Documents/GitHub/ml-etivity3$ ./run.sh
Attivazione dell'ambiente virtuale.
./run.sh: line 26: venv-etivity3/bin/activate: No such file or directory
Aggiornamento di pip...
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: pip in /home/simone/.local/lib/python3.10/site-packages (25.1)
Installazione delle dipendenze da requirements.txt...
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: pandas=1.3.3 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 1)) (1.3.5)
Requirement already satisfied: matplotlib=3.4.3 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 2)) (3.4.3)
Requirement already satisfied: numpy=1.21.0 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 3)) (1.22.4)
Requirement already satisfied: scikit-learn=1.0 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 4)) (1.6.1)
Requirement already satisfied: seaborn=0.11.2 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 5)) (0.11.2)
Requirement already satisfied: scipy=1.8.0 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 6)) (1.7.3)
Requirement already satisfied: PyQt5=5.15.6 in /home/simone/.local/lib/python3.10/site-packages (from -r requirements.txt (line 7)) (5.15.6)
Requirement already satisfied: PyQT5=5.15.6 in /home/simone/.local/lib/python3.10/site-packages (from PyQT5=5.15.6--r requirements.txt (line 7)) (5.15.6)
Requirement already satisfied: python-dateutil=2.7.3 in /home/simone/.local/lib/python3.10/site-packages (from pandas=1.3.3--r requirements.txt (line 1)) (2.9.0.post0)
Requirement already satisfied: pytz=2017.3 in /usr/lib/python3/dist-packages (from pandas=1.3.3--r requirements.txt (line 1)) (2022.4)
Requirement already satisfied: cycler=0.10 in /home/simone/.local/lib/python3.10/site-packages (from matplotlib=3.4.3--r requirements.txt (line 4)) (0.12.1)
Requirement already satisfied: kiwisolver=1.0.1 in /home/simone/.local/lib/python3.10/site-packages (from matplotlib=3.4.3--r requirements.txt (line 4)) (1.4.8)
Requirement already satisfied: pillow=6.2.0 in /usr/lib/python3/dist-packages (from matplotlib=3.4.3--r requirements.txt (line 4)) (9.0.1)
Requirement already satisfied: pyparsing=2.2.1 in /usr/lib/python3/dist-packages (from matplotlib=3.4.3--r requirements.txt (line 4)) (2.4.7)
Requirement already satisfied: joblib=1.2.0 in /home/simone/.local/lib/python3.10/site-packages (from scikit-learn=1.0--r requirements.txt (line 3)) (1.4.2)
Requirement already satisfied: threadpoolctl=3.1.0 in /home/simone/.local/lib/python3.10/site-packages (from scikit-learn=1.0--r requirements.txt (line 3)) (3.6.0)
Requirement already satisfied: six=1.5 in /usr/lib/python3/dist-packages (from python-dateutil=2.7.3--pandas=1.3.3--r requirements.txt (line 1)) (1.16.0)
Avvio del programma...

ML-Etivity3

Tipi di variabili rilevati:
- sepal length (cm): numeric_continuous
- sepal width (cm): numeric_continuous
- petal length (cm): numeric_continuous
- petal width (cm): numeric_continuous
- target: categorical_nominal

=== ANALISI AUTOMATICA IN CORSO ===

1. ANALISI ESPLORATIVA:

Statistiche descrittive:
sepal length (cm) sepal width (cm) petal length (cm) petal width (cm) target
count 150.000000 150.000000 150.000000 150.000000 150
unique NaN NaN NaN NaN 3
top NaN NaN NaN NaN setosa
freq NaN NaN NaN NaN 50
mean 5.843333 3.877333 3.758000 1.199333 NaN
std 0.828066 0.435866 1.765298 0.762238 NaN
min 4.300000 2.000000 1.000000 0.100000 NaN
25% 5.100000 2.800000 1.600000 0.300000 NaN
50% 5.800000 3.800000 4.350000 1.300000 NaN
75% 6.400000 3.300000 5.100000 1.800000 NaN
max 7.900000 4.400000 6.900000 2.500000 NaN

Matrice di correlazione:
sepal length (cm) sepal width (cm) petal length (cm) petal width (cm)
sepal length (cm) 1.000000 -0.117570 0.871754 0.817941
sepal width (cm) -0.117570 1.000000 -0.428440 -0.366126
petal length (cm) 0.871754 -0.428440 1.000000 0.962865
petal width (cm) 0.817941 -0.366126 0.962865 1.000000
QCoreApplication::exec: The event loop is already running

Distribuzione di frequenza per target:
setosa 50
versicolor 50
virginica 50
Name: target, dtype: int64
QCoreApplication::exec: The event loop is already running

2. ANALISI NON SUPERVISIONATA:

[1/2] CLUSTERING K-MEANS
- k=2: Silhouette=0.582
- k=3: Silhouette=0.469
- k=4: Silhouette=0.387
- k=5: Silhouette=0.346
- k=6: Silhouette=0.317
- k=7: Silhouette=0.320
- k=8: Silhouette=0.339
- k=9: Silhouette=0.342
- k=10: Silhouette=0.352

Numero ottimale di cluster: 2 (silhouette=0.582)

[2/2] CLUSTERING GERARCHICO
Taglio ottimale (con offset=2) a distanza: 8.01 (3 cluster)
QCoreApplication::exec: The event loop is already running
QCoreApplication::exec: The event loop is already running

=== ANALISI COMPLETATA ===

```

Figure 5.1: Esecuzione su terminale Linux

5.1 Flusso Operativo

Il processo di analisi prevede una sequenza logica di operazioni:

5.1.1 Selezione del Dataset

L'utente sceglie tra i dataset predefiniti di scikit-learn (Iris, Wine, Breast Cancer, Diabetes o California Housing) attraverso un menu a tendina intuitivo.

5.1.2 Configurazione dell'Analisi

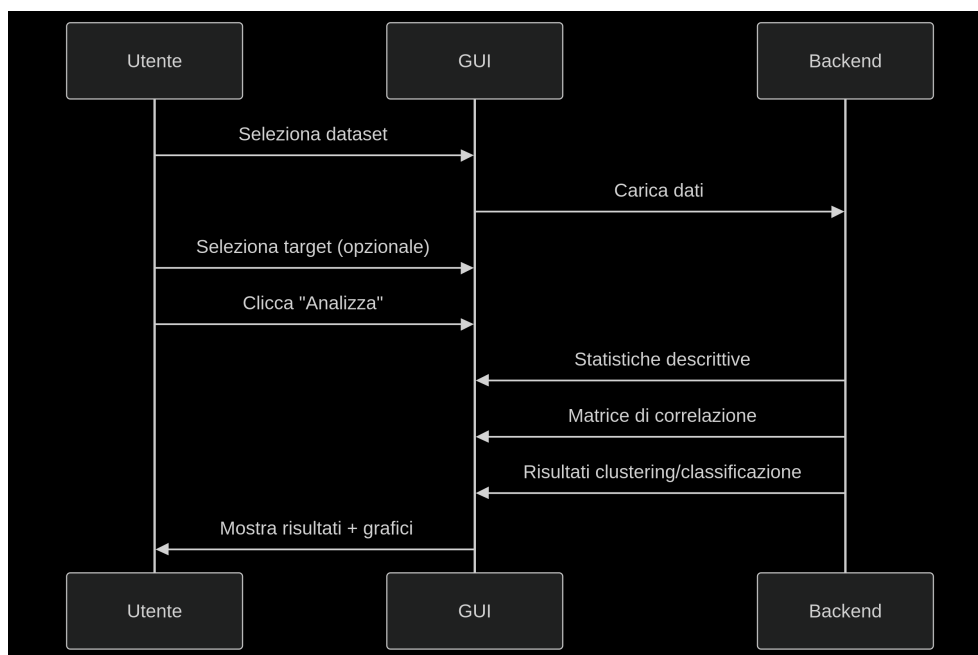
- Scelta opzionale della variabile target per analisi supervisionate
- Possibilità di analisi non supervisionata lasciando il campo target vuoto

5.1.3 Esecuzione dell'Analisi

- Automatica classificazione delle variabili
- Preprocessing dei dati (codifica e normalizzazione)
- Applicazione degli algoritmi appropriati in base alla configurazione

5.1.4 Visualizzazione dei Risultati

- Output testuale nell'area dedicata e grafici mostrati a schermo.



5.2 Interazione Avanzata

Per analisi specifiche, l'utente può:

- Utilizzare il pulsante *"Analizza relazione tra due variabili"* per test di dipendenza (Test Chi-Quadro)
- Visualizzare dendrogrammi per il clustering gerarchico
- Esaminare le feature importance nei modelli supervisionati

5.3 Output e Risultati

Il sistema genera:

- Report statistici completi
- Visualizzazioni heatmap, scatter plot, bar chart
- Metriche di valutazione per i modelli ML
- Risultati del test Chi-quadro con interpretazione statistica

Tutti i risultati rimangono disponibili nell'interfaccia fino alla chiusura dell'applicazione, permettendo all'utente di analizzare comodamente i diversi aspetti dei dati.

