



CAR.CSV

Analisi Statistica

CONTENTS



1

Presentazione del dataset

2

Pre-processing e splitting

3

Fase EDA

4

Addestramento e Hyperparameter tuning

5

Valutazione delle performance e studio statistico

		Car						
1		Buying_Price	Maintenance_Price	No_of_Doors	Person_Capacity	Size_of_Luggage	Safety	Car_Acceptability
2		vhigh	vhigh	2	2	small	low	unacc
3		vhigh	vhigh	2	2	small	med	unacc
4		vhigh	vhigh	2	2	small	high	unacc
5		vhigh	vhigh	2	2	med	low	unacc
6		vhigh	vhigh	2	2	med	med	unacc
7		vhigh	vhigh	2	2	med	high	unacc
8		vhigh	vhigh	2	2	big	low	unacc
370		vhigh	low	3	4	big	high	acc
371		vhigh	low	3	more	small	low	unacc
372		vhigh	low	3	more	small	med	unacc
373		vhigh	low	3	more	small	high	acc
374		vhigh	low	3	more	med	low	unacc
375		vhigh	low	3	more	med	med	acc
376		vhigh	low	3	more	med	high	acc
377		vhigh	low	3	more	big	low	unacc
378		vhigh	low	3	more	big	med	acc

[5 rows x 7 columns]

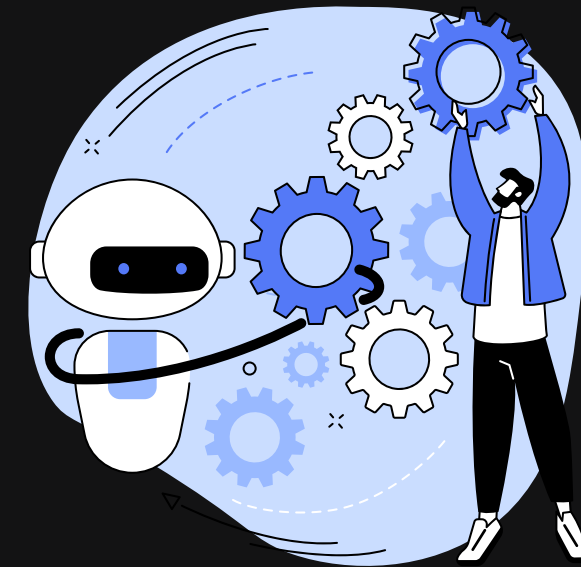
Buying_Price	0
Maintenance_Price	0
No_of_Doors	0
Person_Capacity	0
Size_of_Luggage	0
Safety	0
Car_Acceptability	0



Presentazione del dataset, composizione:

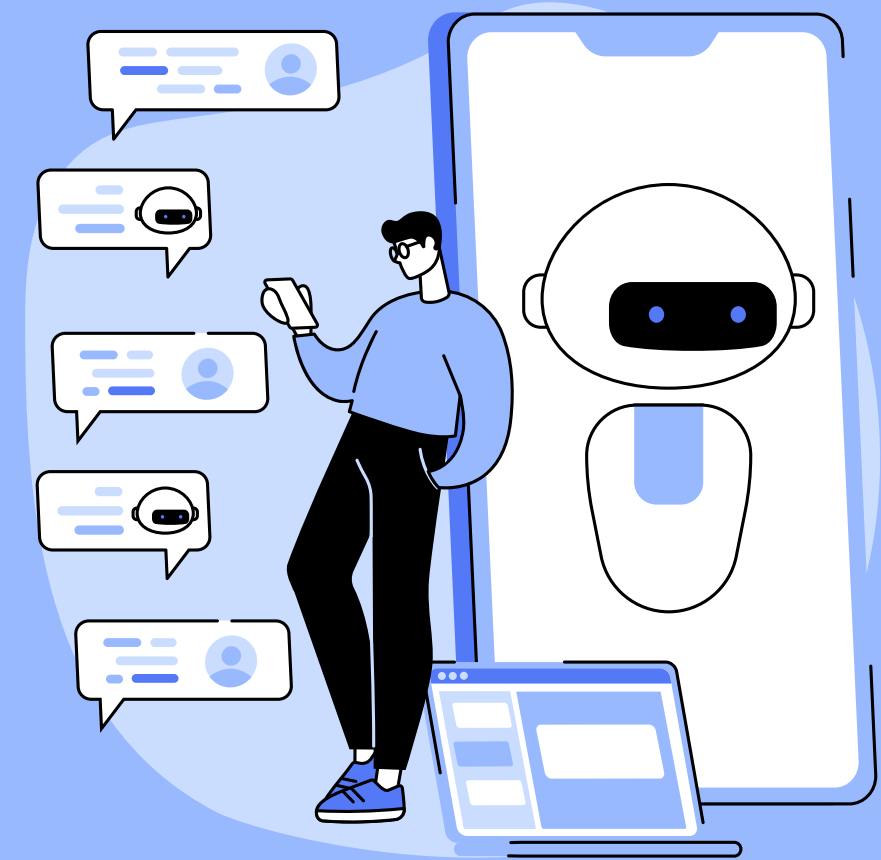
Fase PRE-PROCESSING:

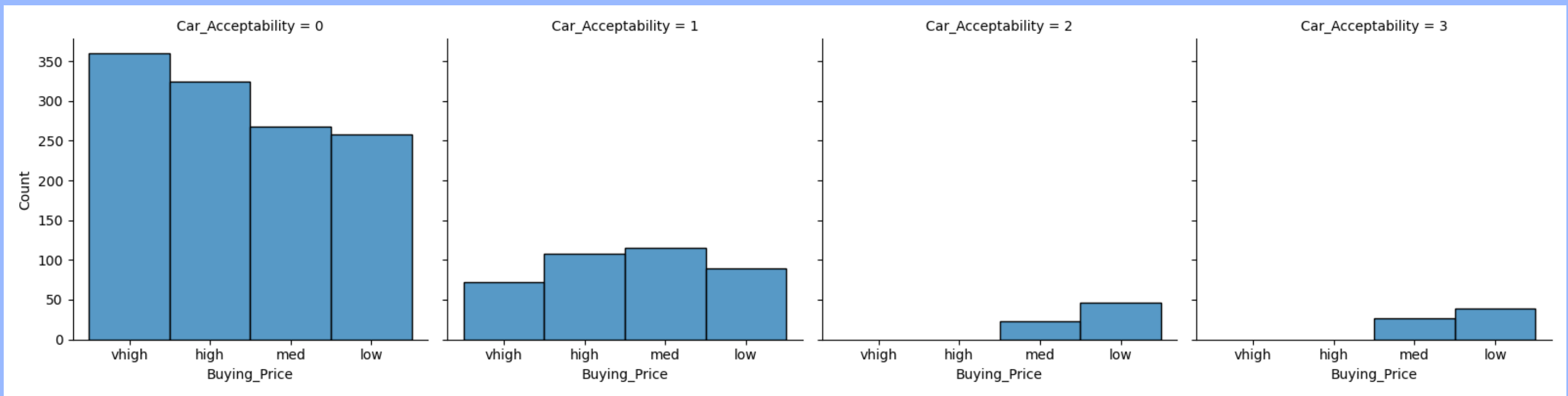
```
33 print(dataframe.shape)
34 print(dataframe.head())
35 print(dataframe.isnull().sum())
36
37
38
39 # tutto ottimo, procediamo con il controllo dei valori fuori soglia
40
41 min_limit = 0 # Valore minimo accettabile generale, potremmo farne uno specifico ma teniamoci indicativi
42 max_limit = 5 # Valore massimo accettabile per un dataset sulle auto
43
44 # Seleziona solo le colonne numeriche
45 numeric_columns = ['No_of_Doors', 'Person_Capacity']
46
47 # Itera sulle colonne numeriche e controlla i valori fuori soglia
48 for column in numeric_columns:
49     if column == 'No_of_Doors':
50         outliers = dataframe[~dataframe[column].str.contains('5more')]
51     else:
52         outliers = dataframe[(dataframe[column] < min_limit) | (dataframe[column] > max_limit)]
53         if not outliers.empty:
54             print(f"Variabile: {column}")
55             print(outliers)
56             print("\n")
57
```



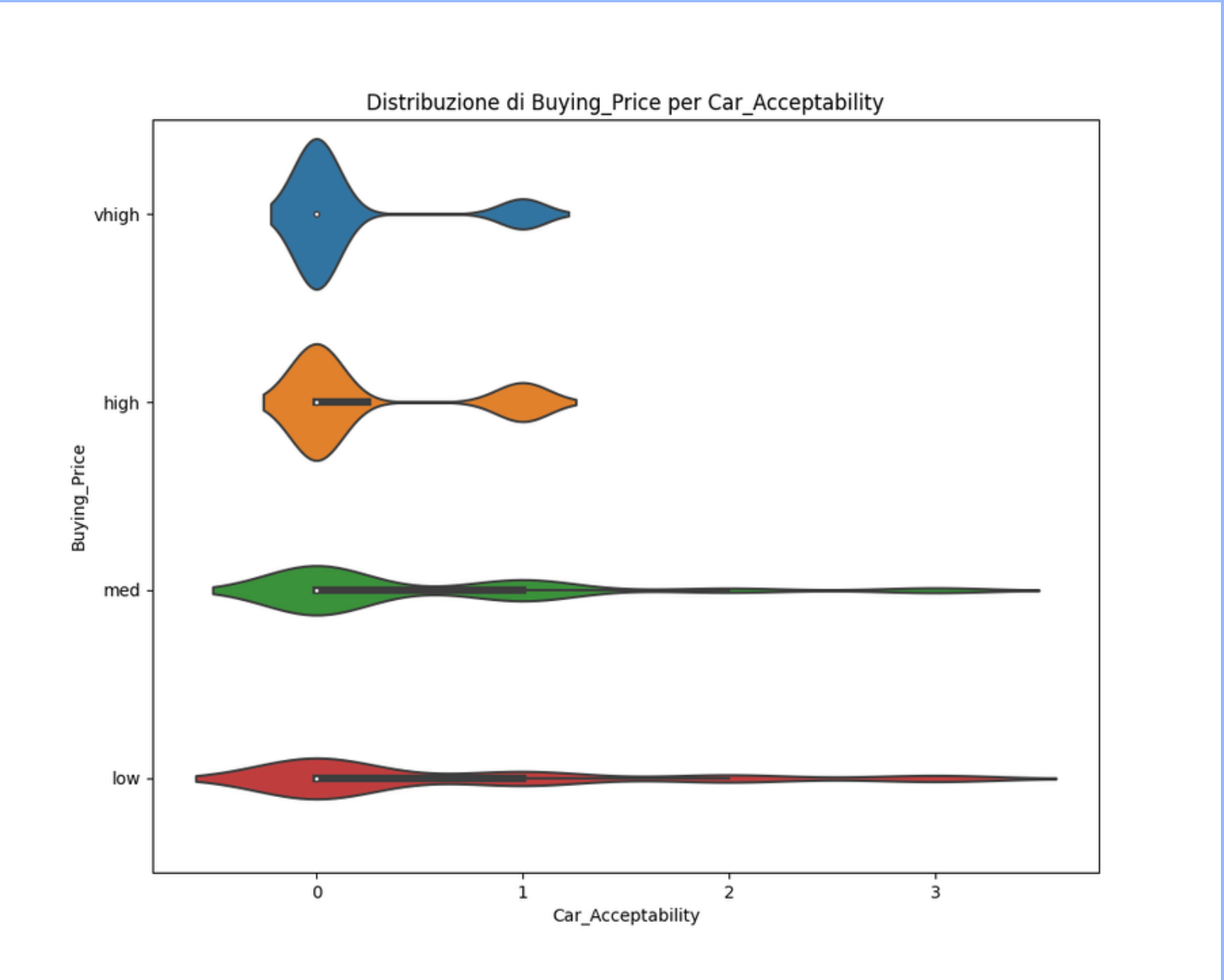
Fase di SPLITTING:

```
132 # tutto andato a buon fine, ma prima di procedere con lo splitting
133 print(dataframe.head())
134
135 X = dataframe.drop(columns='Car_Acceptability')
136 y = dataframe['Car_Acceptability']
137 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
138
```





Possiamo notare con una rapida analisi EDA con i boxplot che al diminuire della possibilità di acquisto da parte di un utente finale corrisponda un'accettabilità del veicolo stesso più elevata, una visione più dinamica può esser vista nel seguente grafico interattivo: [localhostlink](#)



FASE EDA:

Addestramento dei modelli e confronto dei due algoritmi:

```
139 #modello regressione logistica
140 modelloLogistica = LogisticRegression(solver='lbfgs', max_iter=1000)
141 modelloLogistica.fit(X_train, y_train)
142 # Definisci il modello SVC
143 model = SVC()
144 model.fit(X_train, y_train)
145
146 # Effettua predizioni sui dati di test
147 y_pred_logistic = modelloLogistica.predict(X_test)
148 #Adesso predizioni Svc
149 y_pred_svc = model.predict(X_test)
150 # Valuta le prestazioni del modello
151 accuracy_logistic = accuracy_score(y_test, y_pred_logistic)
152 print("Accuracy del modello di regressione logistica:", accuracy_logistic)
153 classification_report_logistic = classification_report(y_test, y_pred_logistic)
154 print("Report di classificazione del modello di regressione logistica:")
155 print(classification_report_logistic)
156 #ADESSO PARTE SVC
157 accuracy_SVC = accuracy_score(y_test, y_pred_svc)
158 print("Accuracy del modello SVC:", accuracy_SVC)
159 classification_report_SVC = classification_report(y_test, y_pred_svc)
160 print("Report di classificazione del modello SVC:")
161
```

[5 rows x 7 columns]

Accuracy del modello di regressione logistica: 0.8458574181117534

Report di classificazione del modello di regressione logistica:

	precision	recall	f1-score	support
0	0.90	0.94	0.92	358
1	0.70	0.64	0.67	118
2	0.62	0.42	0.50	19
3	0.83	0.79	0.81	24
accuracy			0.85	519
macro avg	0.76	0.70	0.72	519
weighted avg	0.84	0.85	0.84	519

Accuracy del modello SVC: 0.928709055876686

Report di classificazione del modello SVC:

Media: 0.8895

Mediana: 0.8786

Minimo: 0.8266

Massimo: 0.9827

Intervallo: 0.1561

Deviazione standard: 0.0435

Hyperparameter tuning:

```
164 #Addestramento del modello SVM con kernel rbf
165 param_grid_rbf = {
166     'C': [0.1, 1.0, 10.0],
167     'gamma': ['scale', 'auto']
168 }
169 grid_search_rbf = GridSearchCV(estimator=SVC(kernel='rbf'), param_grid=param_grid_rbf, cv=5)
170 grid_search_rbf.fit(X_train, y_train)
171 best_model_rbf = grid_search_rbf.best_estimator_
172 accuracy_rbf = best_model_rbf.score(X_test, y_test)
173
174 # Addestramento del modello SVM con kernel polinomiale
175 param_grid_poly = {
176     'C': [0.1, 1.0, 10.0],
177     'gamma': ['scale', 'auto'],
178     'degree': [2, 3, 4]
179 }
180 grid_search_poly = GridSearchCV(estimator=SVC(kernel='poly'), param_grid=param_grid_poly, cv=5)
181 grid_search_poly.fit(X_train, y_train)
182 best_model_poly = grid_search_poly.best_estimator_
183 accuracy_poly = best_model_poly.score(X_test, y_test)
184
185 # Addestramento del modello SVM con kernel lineare
186 param_grid_linear = {
187     'C': [0.1, 1.0, 10.0]
188 }
189 grid_search_linear = GridSearchCV(estimator=SVC(kernel='linear'), param_grid=param_grid_linear, cv=5)
190 grid_search_linear.fit(X_train, y_train)
191 best_model_linear = grid_search_linear.best_estimator_
192 accuracy_linear = best_model_linear.score(X_test, y_test)
193
```

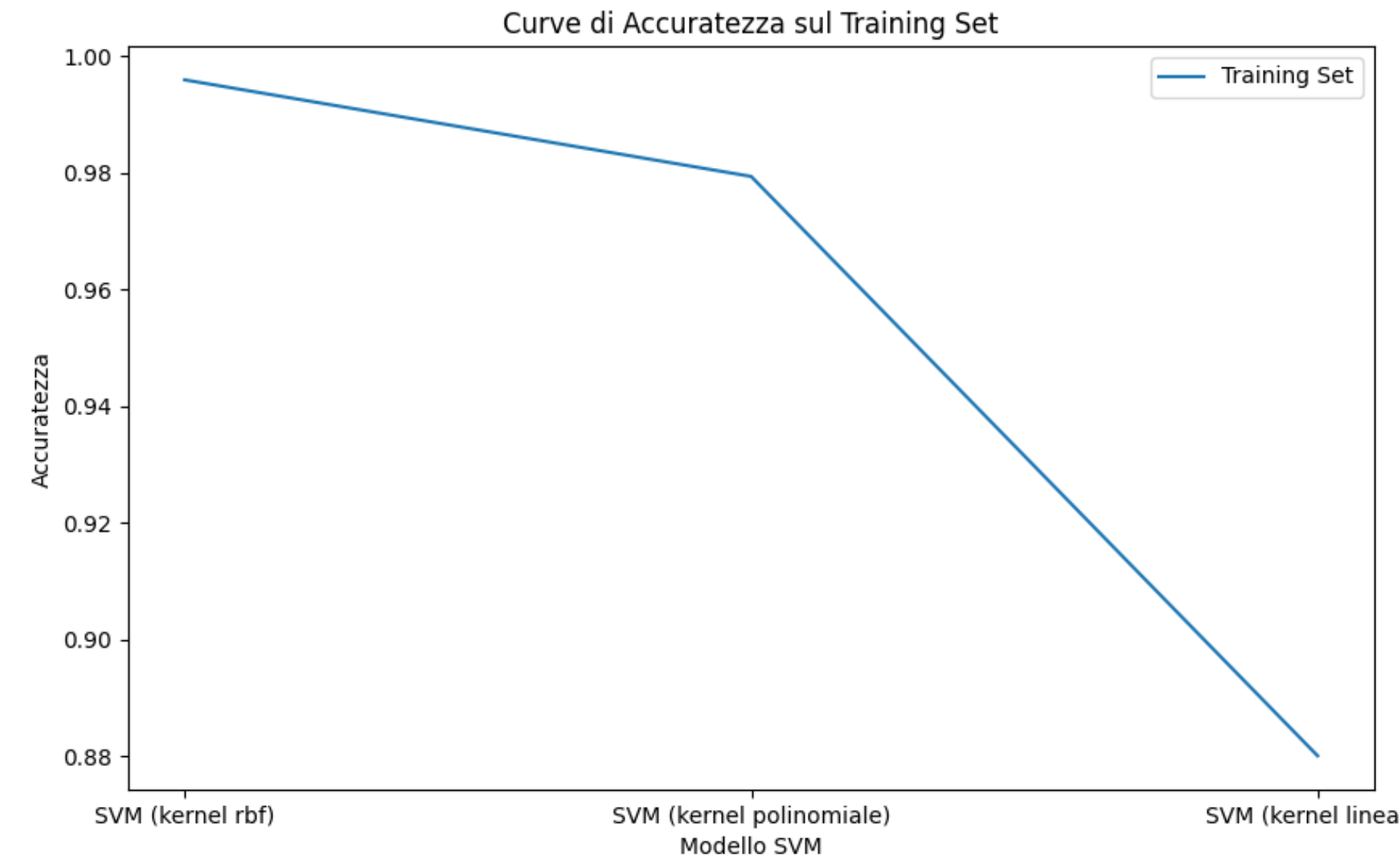
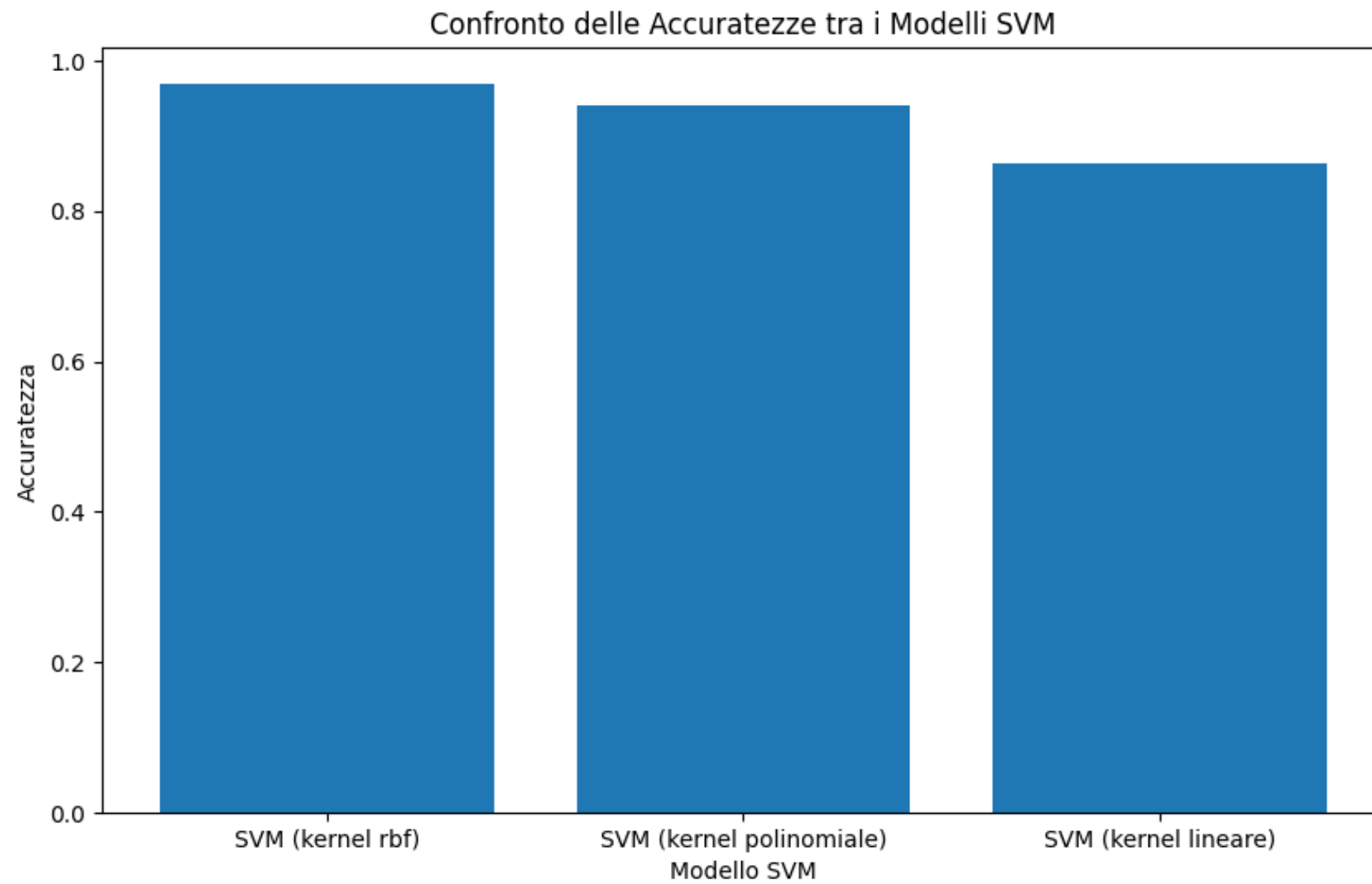


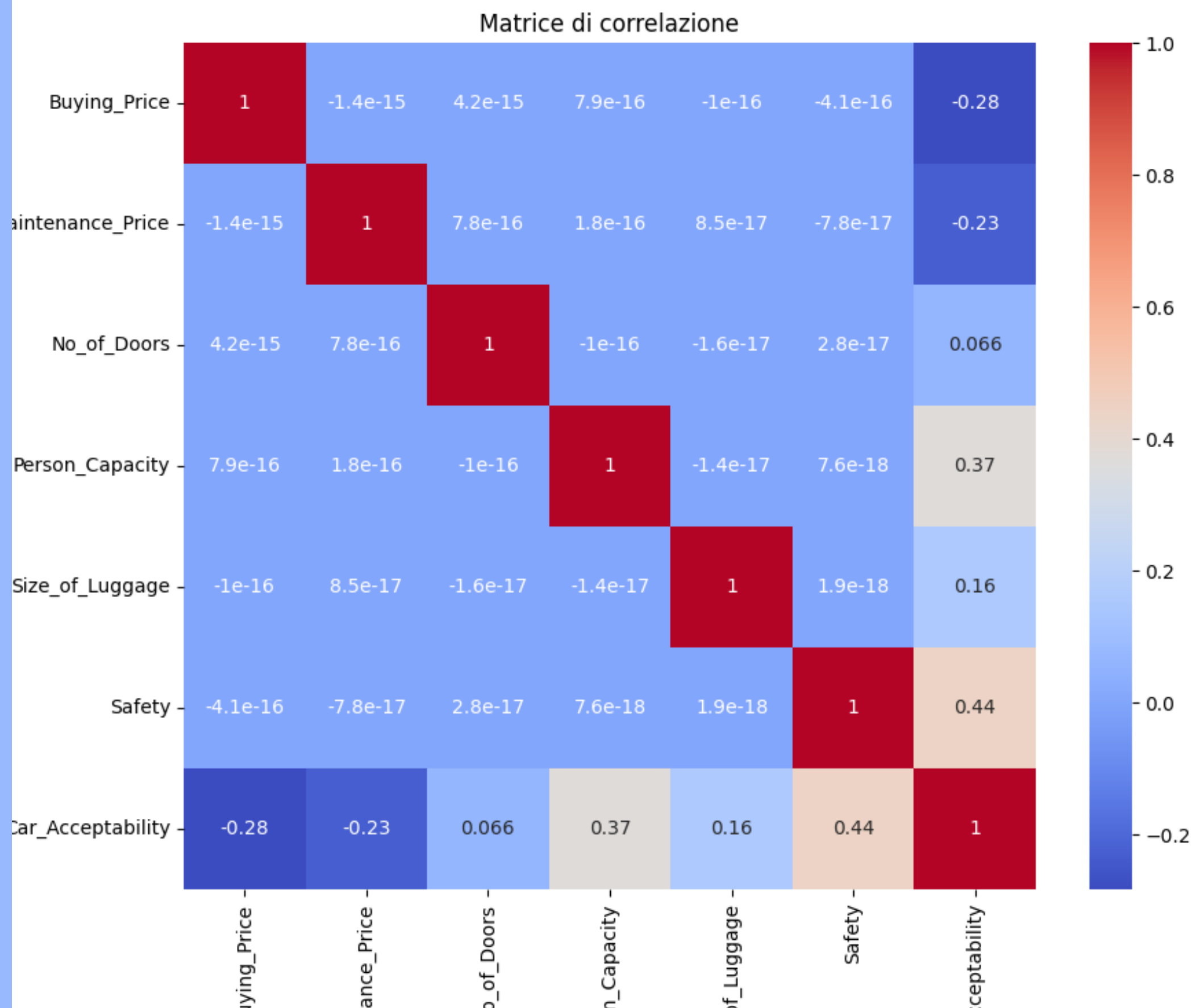
```

208 plt.bar(labels, accuracies)
209 plt.xlabel('Modello SVM')
210 plt.ylabel('Accuratezza')
211 plt.title('Confronto delle Accuratezze tra i Modelli SVM')
212 plt.show()
213
214 # Calcolo delle accuratezze dei modelli sul training set
215 accuracy_rbf_train = best_model_rbf.score(X_train, y_train)
216 accuracy_poly_train = best_model_poly.score(X_train, y_train)
217 accuracy_linear_train = best_model_linear.score(X_train, y_train)
218
219 # Creazione dei vettori di accuratezza
220 accuracies_train = [accuracy_rbf_train, accuracy_poly_train, accuracy_linear_train]
221
222
223 # Etichette per i modelli
224 labels = ['SVM (kernel rbf)', 'SVM (kernel polinomiale)', 'SVM (kernel lineare)']
225
226 # Grafico delle curve di accuratezza per il training set
227 plt.figure(figsize=(10, 6))
228 plt.plot(accuracies_train, label='Training Set')
229 plt.xticks(range(len(labels)), labels)
230 plt.xlabel('Modello SVM')
231 plt.ylabel('Accuratezza')
232 plt.title('Curve di Accuratezza sul Training Set')
233 plt.legend()
234 plt.show()
235
236
237 # Esegui la validazione incrociata con k ripetizioni (k >= 10)
238 k = 10
239 scores = cross_val_score(model, X, y, cv=k, scoring='accuracy')
240

```

Plot in cui risulta evidente una maggiore accuratezza del modello kernel rbf per il nostro dataset.

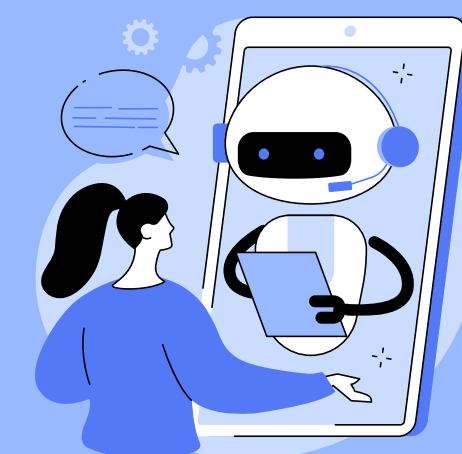




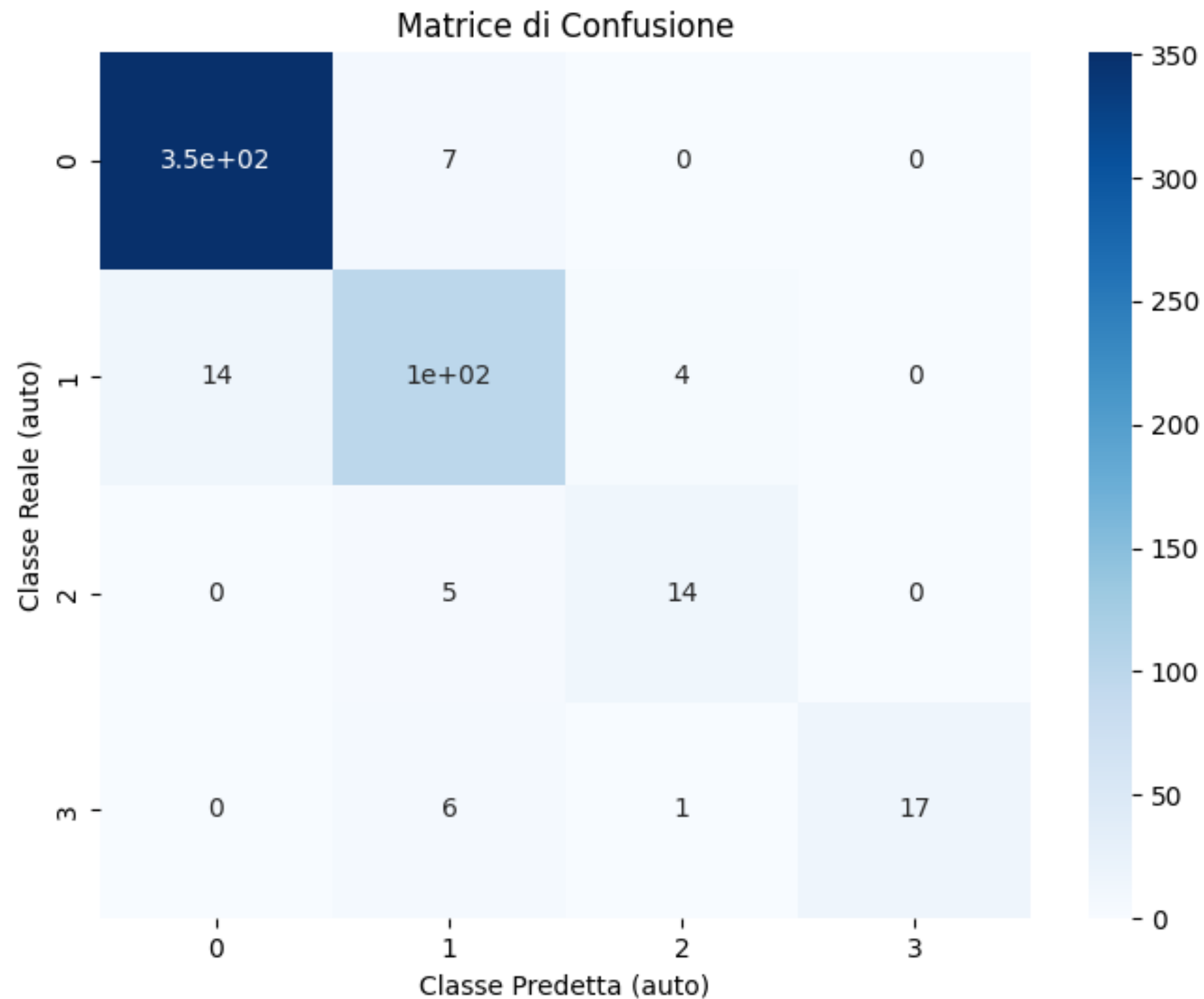
La matrice di correlazione ci consente di individuare le variabili che sono correlate tra loro. Una correlazione positiva significa che le variabili aumentano o diminuiscono insieme, mentre una correlazione negativa significa che le variabili si muovono in direzioni opposte.

Notiamo che il colore rosso indica correlazione positiva perfetta mentre blu scuro quella negativa perfetta, ovvero rispettivamente 1 e -1

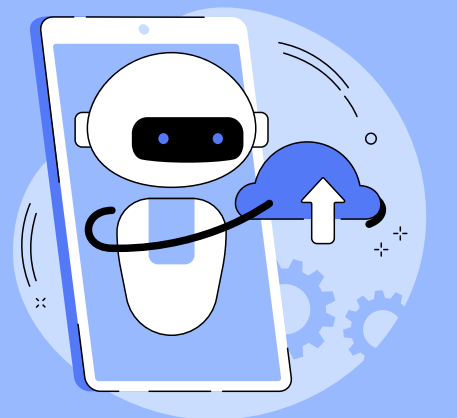
la correlazione non implica causalità. È importante considerare il contesto del tuo dataset e valutare la relazione tra le variabili in base alla conoscenza di dominio e agli obiettivi dell'analisi.



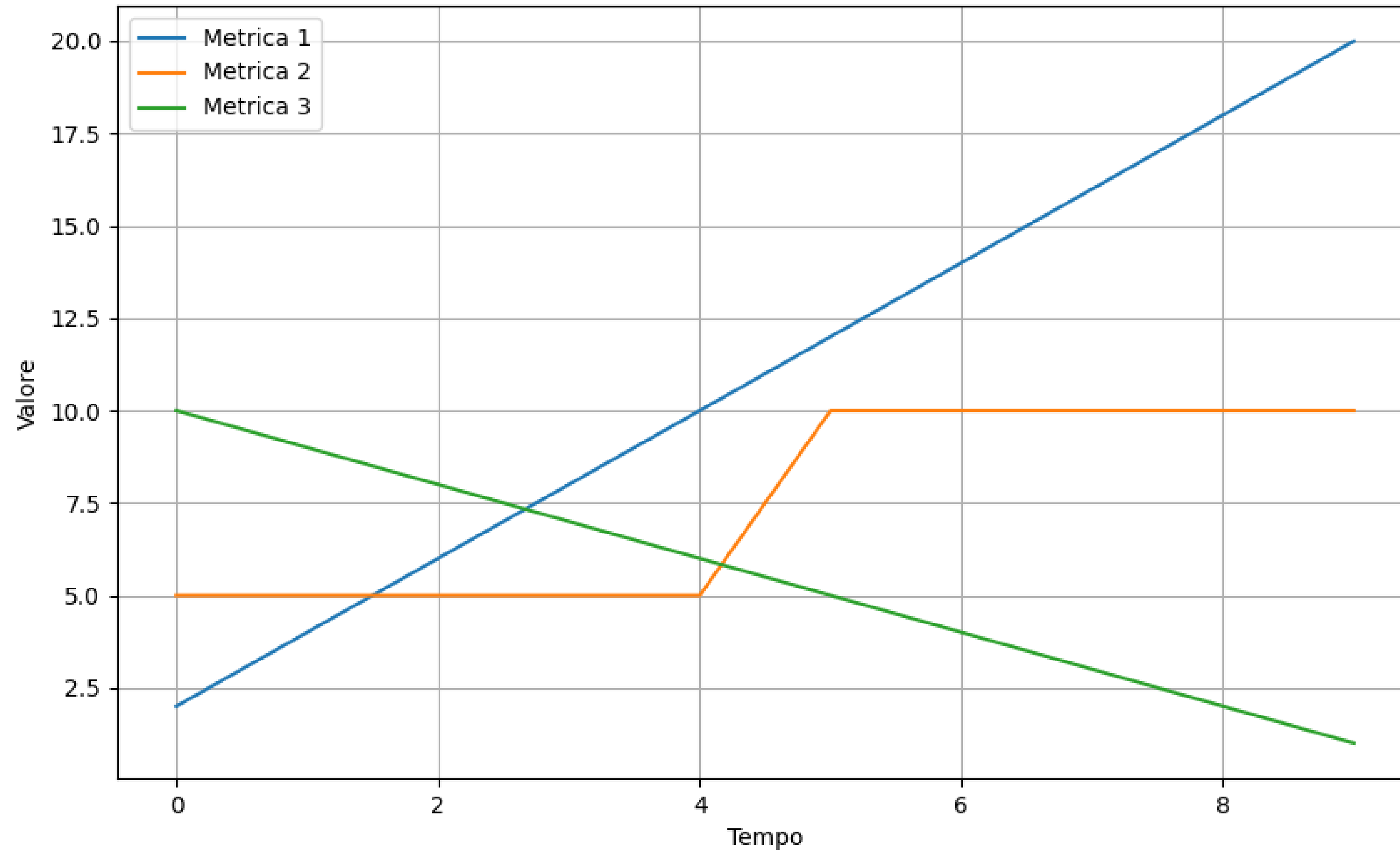
Valutazione delle performance con matrice di confusione.



Come notiamo abbiamo in posizione 0-0 ben 351 istanze che appartengono alla classe reale "unacc" e sono state predette correttamente come "unacc" (True Negatives, TN). Ci sono anche 7 istanze di classe "unacc" che sono state predette erroneamente come altre classi. Non ci sono istanze di altre classi che sono state predette come "unacc" (valori 0 nella riga).

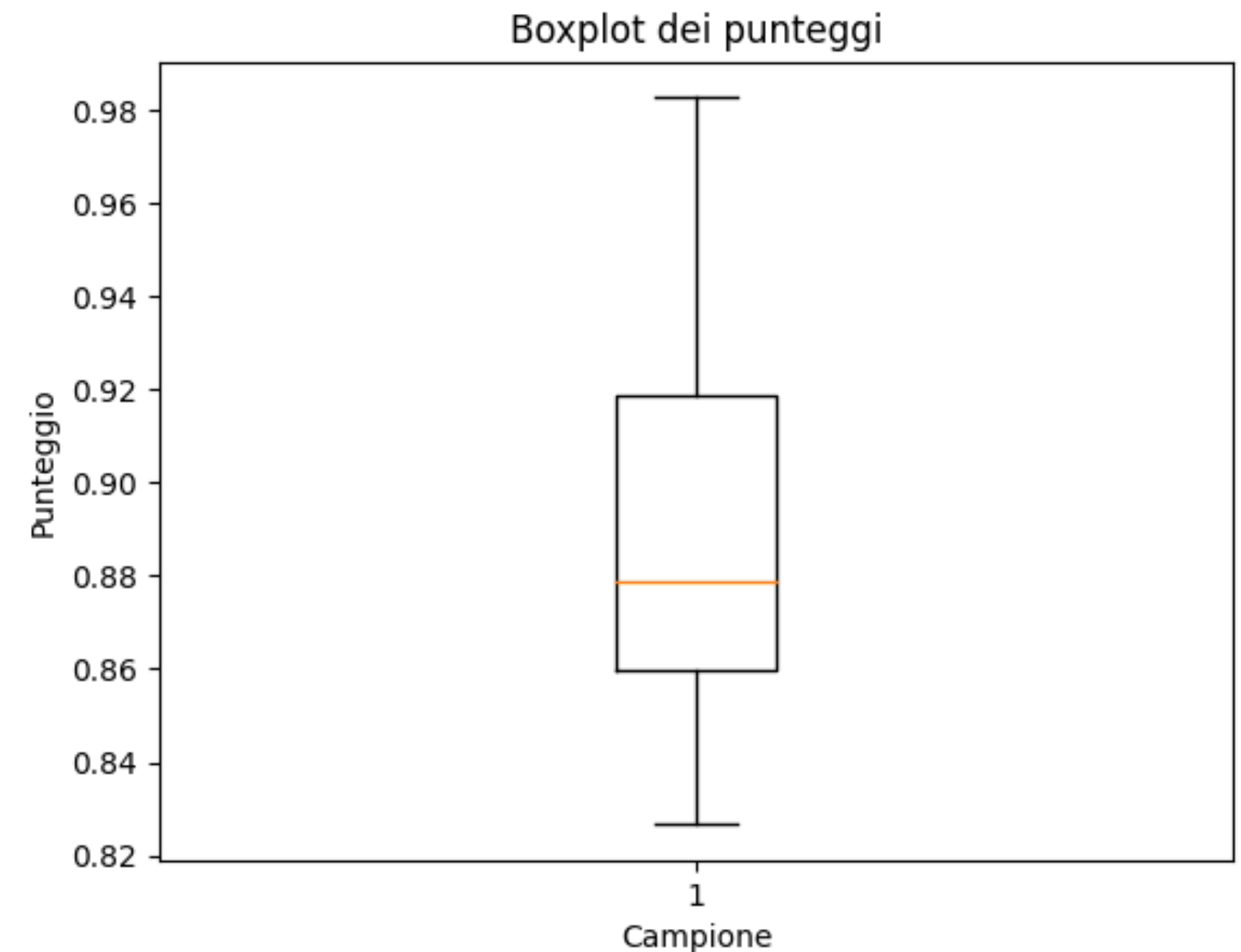
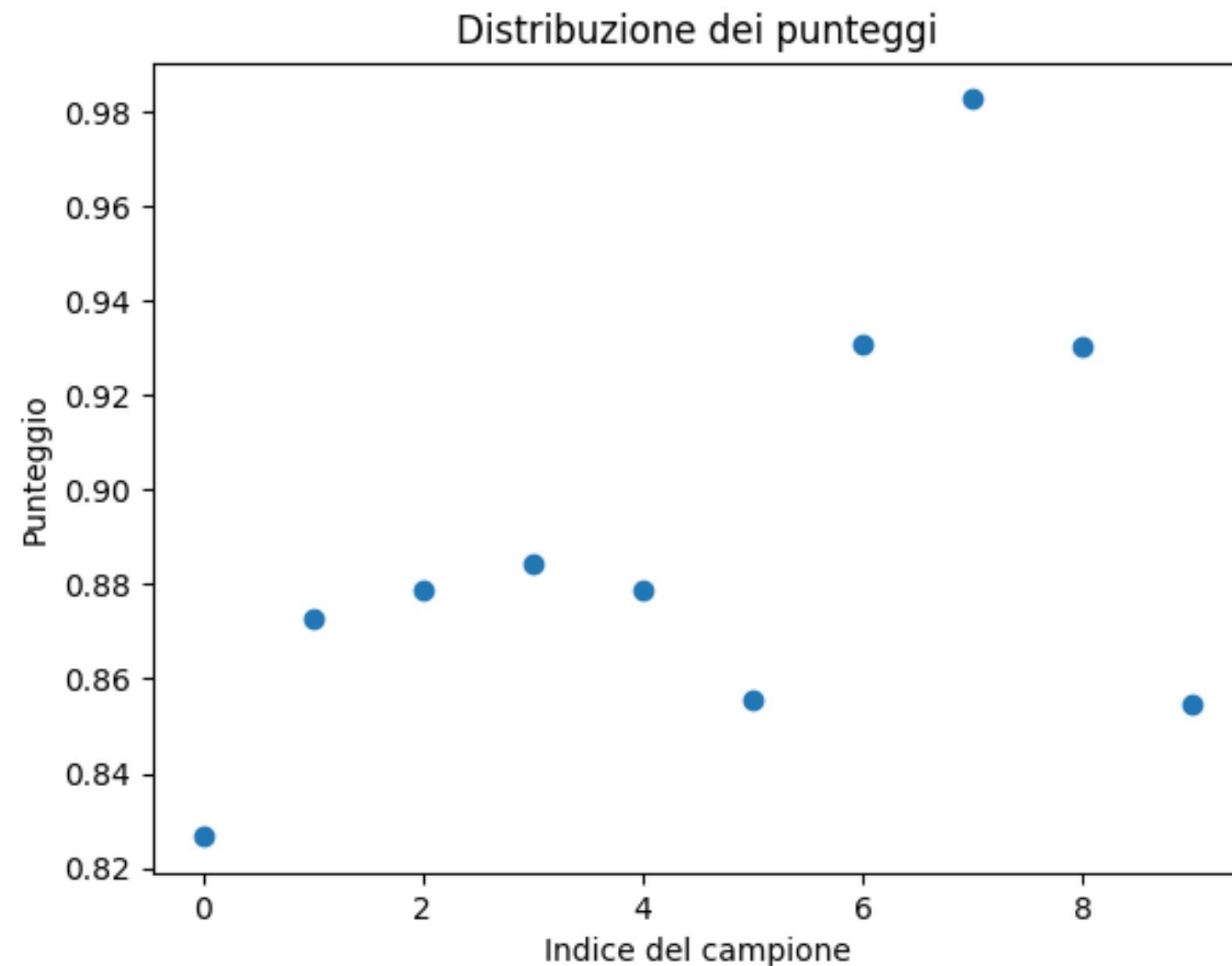


Andamento delle metriche nel tempo

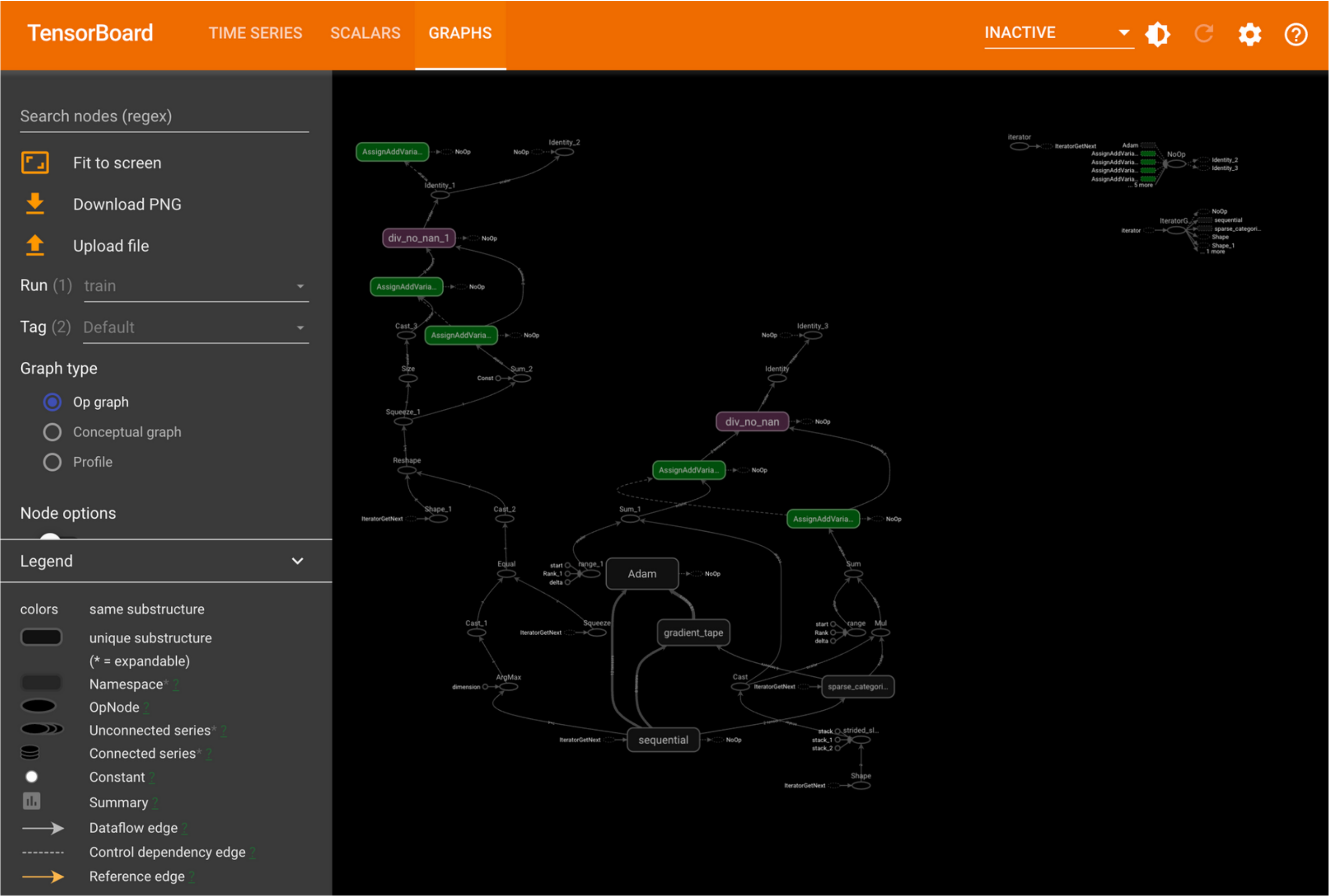


Studio statistico sui risultati della valutazione:

Media: 0.8895
Mediana: 0.8786
Minimo: 0.8266
Massimo: 0.9827
Intervallo: 0.1561
Deviazione standard: 0.0435
Intervallo di confidenza al 95%: [0.8584, 0.9206]



Applicazione di Rete Neurale con Tensorflow:





**PROGETTO DI MASSIMO ANDRES ROJAS
CAMACHO.**

NUMERO MATRICOLA: 0001019759

MASSIMO.ROJASCAMACHO@STUDIO.UNIBO.IT