

# Progetto per ingegneria della conoscenza

Componenti del gruppo:

- Ciavarella Andrea, [760411], [a.ciavarella16](#)

Link repository github: <https://github.com/AndreCiav/icon.git>

A.A 2023-2024



## INDICE

- 1. Introduzione**
  - 1.1. Strumenti utilizzati
- 2. Dataset e preprocessing**
  - 2.1. Preprocessing del dataset
- 3. Apprendimento supervisionato**
  - 3.1. Scelta dei modelli
  - 3.2. Scelta degli iper-parametri
  - 3.3. Valutazione modelli
  - 3.4. Varianza e deviazione standard
- 4. Base di conoscenza e Prolog**
  - 4.1. Conversione da albero a Prolog
  - 4.2. Regole e fatti
- 5. Apprendimento bayesiano**
  - 5.1. Struttura della rete bayesiana
  - 5.2. Addestramento della rete bayesiana
  - 5.3. Utilizzo della rete bayesiana
- 6. Conclusioni**
  - 6.1. Possibili sviluppi
  - 6.2. Riferimenti bibliografici

## 1. INTRODUZIONE

Il dataset considerato per questo progetto si può trovare al link <https://www.kaggle.com/datasets/kunwarakash/chennai-housing-sales-price>, contiene dati di vendita di proprietà immobiliari, e tramite queste caratteristiche si vuole costruire un modello di apprendimento automatico in grado di prevedere il prezzo di vendita delle case e costruire una base di conoscenza complessiva. Successivamente verranno indicate e descritte tutte le caratteristiche del dataset.

### 1.1 Strumenti utilizzati

Il progetto è stato sviluppato in Python, linguaggio di programmazione che offre una serie di librerie utili e adatte allo studio, all'analisi e apprendimento di enormi quantità di dati. L'ambiente di sviluppo usato è Pycharm, Visual Studio Code e Colab. Librerie python utilizzate:

- **Pandas**: usato per organizzare e trasformare i dati in strutture, permettendo di filtrare, pulire e preparare i dati importati dal file CSV per l'analisi.
- **Scikit-learn** (sklearn): usato per l'apprendimento automatico, offre una vasta gamma di algoritmi di apprendimento supervisionato e non supervisionato, oltre a strumenti per la selezione e la valutazione dei modelli.

- **Seaborn:** aiuta a creare visualizzazioni statistiche, come le heatmap e i pairplot, per osservare le relazioni tra le variabili nel dataset.
- **Matplotlib:** Costruisce grafici personalizzati, come le curve di apprendimento o le distribuzioni, per monitorare le performance dei modelli e interpretare i dati.

## 2. DATASET E PREPROCESSING

Il dataset è formato dai seguenti campi(22 colonne,7109 righe):

- **PRT\_ID:** Identificativo unico della proprietà.
- **AREA:** Area o località(paese) in cui si trova la proprietà.
- **INT\_SQFT:** Superficie interna in piedi quadrati.
- **DATE\_SALE:** Data di vendita della proprietà.
- **DIST\_MAINROAD:** Distanza della proprietà dalla strada principale in metri.
- **N\_BEDROOM:** Numero di camere da letto.
- **N\_BATHROOM:** Numero di bagni.
- **N\_ROOM:** Numero totale di stanze.
- **SALE\_COND:** Condizione di vendita (ad es. nuovo, rivendita).
- **PARK\_FACIL:** Indica se è presente un'area di parcheggio.
- **DATE\_BUILD:** Data di costruzione della proprietà.
- **BUILDTYPE:** Tipo di costruzione (ad es. casa indipendente, appartamento).
- **UTILITY\_AVAIL:** Servizi disponibili (es. acqua, elettricità).
- **STREET:** Tipo di strada (ad es. strada asfaltata, ghiaia).
- **MZZONE:** Zona in cui si trova la proprietà.
- **QS\_ROOMS:** Qualità delle stanze in valori da (1-10).
- **QS\_BATHROOM:** Qualità del bagno in valori da (1-10).
- **QS\_BEDROOM:** Qualità della camera da letto in valori da (1-10).
- **QS\_OVERALL:** Qualità complessiva in valori da (1-10).
- **REG\_FEE:** Importo della tassa di registrazione associata alla transazione di vendita della proprietà, può variare in base al valore della proprietà.
- **COMMIS:** Commissione di vendita pagata agli agenti o intermediari.
- **SALES\_PRICE:** Prezzo di vendita della proprietà.

### 2.1 Pre-processing del dataset

L'obiettivo del preprocessing del dataset è migliorare la qualità dei dati per renderli adatti all'analisi e alla modellazione. Questo processo include la pulizia dei dati, la trasformazione delle variabili, la gestione degli outlier e la selezione delle caratteristiche più rilevanti.

Alcuni valori all'interno del database si presentano in modo errato, a causa di errori di battitura, quindi vengono effettuate le varie correzioni/rinominazioni.

```

'SALE_COND': {
  # Unifica le variazioni di scrittura per le condizioni di vendita
  'AbNormal': 'Ab Normal',
  'Ab Normal': 'Ab Normal',
  'Partiall': 'Partial',
  'Adj Land': 'AdjLand',
  'PartiaLl': 'Partial'
},
'PARK_FACIL': {
  # Corregge errori di battitura nella disponibilità di parcheggio
  'Noo': 'No'
},
'BUILDTYPE': {
  # Corregge errori di battitura nei tipi di costruzione
  'Comercial': 'Commercial',
  'Other': 'Others'
},

```

I valori corretti vengono salvati in un file *Chennai\_Sales\_Data\_Corrected.csv*

Le colonne DATE\_SALE (data di vendita) e DATE\_BUILD (data di costruzione) sono convertite in formato datetime, permettendo di estrarre informazioni utili come l'anno e il mese, quindi vengono create due colonne SALE\_YEAR e SALE\_MONTH.

Successivamente vengono identificate le variabili categoriche, ossia variabili che contengono testo (di solito le analisi vengono effettuate con valori numerici, quindi vanno convertiti), e se utili vengono riadattate al database con varie tecniche (label encoders oppure one-hot encoding), se non utili vengono eliminati o non considerate a causa della poca importanza.

La tecnica utilizzata è label encoders ovvero, viene trasformata una colonna con valori categorici in una colonna numerica. Questo metodo assegna un numero intero univoco a ciascun valore unico nella colonna originale. Per esempio, se la colonna "AREA" contiene valori come "Adyar", "Anna Nagar" e "Velachery", il Label Encoding assegnerà numeri come 1, 2, e 3 rispettivamente a ciascuno di questi valori.

Nel nostro caso le variabili categoriche con i rispettivi valori sono:

- 'PRT\_ID': id dell'appartamento quindi , 1,2,3,4,5,6,7 ecc.
- 'AREA= ['Adyar', 'Anna Nagar', 'Chrompet', 'Karapakkam', 'KK Nagar', 'TNagar', 'T Nagar', 'Velachery'], indicano le città
- 'SALE\_COND',=['AbNormal', 'Family', 'Partial', 'AdjLand', 'Normal Sale']

Specifica la condizione della vendita, con i seguenti possibili valori:

- AbNormal (situazione di vendita non ordinaria),
- Family (vendita per motivi familiari),
- Partial (vendita parziale dell'immobile),
- AdjLand (terreno adiacente incluso nella vendita),
- Normal Sale (vendita standard).

- 'PARK\_FACIL'=['Yes', 'No'], yes se c'è disponibilità di parcheggio, no altrimenti
- 'BUILDTYPE',=['Commercial', 'Others', 'House'], descrive il tipo di edificio in cui si trova l'appartamento

- Commercial (uso commerciale),
  - Others (altre tipologie non specificate),
  - House (residenza)
- 'UTILITY\_AVAIL',=['AllPub', 'NoSeWa'], mostra le utilità disponibili nella proprietà:
    - AllPub (tutti i servizi pubblici disponibili),
    - NoSeWa (mancanza di servizi igienico-sanitari).
  - 'STREET',=['Paved' , 'Gravel', 'No Access'] , indica il tipo di strada d'accesso alla proprietà:
    - Paved (asfaltata),
    - Gravel (ghiaiosa),
    - No Access (nessun accesso diretto).
  - 'MZZONE',=['A', 'RH', 'RL', 'I', 'C' , 'RM'], indicano il codice della zona

Alcune di queste colonne non sono importanti al fine della predizione quindi vengono eliminate

```
# Elimina le caratteristiche che sono meno utili per la previsione o che possono portare all'overfitting
house_data = house_data.drop(columns=['PRT_ID', 'SALE_COND', 'STREET', 'MZZONE', 'UTILITY_AVAIL'])
```

Le restanti vengono codificate:

	AREA	INT_SQFT	DIST_MAINROAD	N_BEDROOM	N_BATHROOM	N_ROOM	PARK_FACIL	BUILDTYPE
0	4	1004.0	131.0	1.0	1.0	3.0	1	0
1	1	1986.0	26.0	2.0	1.0	5.0	0	0
2	0	909.0	70.0	1.0	1.0	3.0	1	0
3	7	1855.0	14.0	3.0	2.0	5.0	0	2
4	4	1226.0	84.0	1.0	1.0	3.0	1	2
...	...	...	...	...	...	...	...	...
7104	4	598.0	51.0	1.0	1.0	2.0	0	2
7105	7	1897.0	52.0	3.0	2.0	5.0	1	2
7106	7	1614.0	152.0	2.0	1.0	4.0	0	1
7107	4	787.0	40.0	1.0	1.0	2.0	1	0
7108	7	1896.0	156.0	3.0	2.0	5.0	1	2

7109 rows x 9 columns

- AREA: variabile che indica il nome della città, convertita in valori da 1 a 8, in base al numero che sono stati identificati, esempio la città Adyar corrisponde a 1, 'Anna Nagar' corrisponde a 2 e così via.
- PARK\_FACIL: variabile che indica la presenza di un parcheggio, codificato in true o false(variabile binaria), perché due valori.
- BUILDTYPE: variabile che indica il tipo di edificio codificato con numeri da 1 a 3

I valori mancanti presenti in alcune colonne vengono sostituiti con la media dei valori, della stessa colonna. Al termine viene salvato il database corretto e preprocessato al fine delle previsioni per i modelli(/preprocessed\_house\_data.csv).

Nel progetto è stato deciso di utilizzare `LabelEncoder` anziché `OneHotEncoder` per semplificare il modello e ridurre il numero di caratteristiche risultanti. `LabelEncoder` assegna un valore numerico intero unico a ciascuna categoria, mantenendo una sola colonna per ogni variabile categoriale. `OneHotEncoder`, invece, crea una colonna separata per ciascuna categoria, il che può risultare in un numero elevato di colonne e aumentare la complessità del modello.

### Split dataset

Il dataset viene suddiviso in set di addestramento e di test per valutare l'efficacia del modello. Questa divisione permette di addestrare il modello su una parte dei dati (set di addestramento 80% del dataset) e poi verificarne le prestazioni su dati separati e mai visti prima (set di test 20% del dataset). In questo modo, si può ottenere una stima più accurata delle capacità predittive del modello su dati nuovi e valutare se è in grado di generalizzare bene, evitando problemi come l'overfitting, ovvero quando il modello si adatta troppo ai dati di addestramento e non funziona altrettanto bene su nuovi dati.

In questa fase, viene caricato il dataset pre-elaborato e vengono separate le caratteristiche (X) dalla variabile target (Y), quindi le feature non conterranno la variabile target 'SALES\_PRICE' rappresentante il prezzo di vendita delle case. Il modello, una volta addestrato, restituisce l'importanza di ciascuna caratteristica, evidenziando quali variabili influenzano maggiormente la previsione del prezzo. Solo le caratteristiche con un'importanza sopra una soglia prestabilita (nel nostro caso `threshold = 0.01`) vengono quindi selezionate, riducendo le variabili meno rilevanti e ottimizzando le prestazioni del modello.

Importanza delle caratteristiche:		
	Feature	Importance
12	REG_FEE	0.770348
0	AREA	0.075496
1	INT_SQFT	0.055797
7	BUILDTYPE	0.035141
3	N_BEDROOM	0.015695
4	N_BATHROOM	0.010283
13	COMMIS	0.009797
2	DIST_MAINROAD	0.003198
11	QS_OVERALL	0.002937
8	QS_ROOMS	0.002755
9	QS_BATHROOM	0.002733
6	PARK_FACIL	0.002572
10	QS_BEDROOM	0.002524
14	SALE_YEAR	0.002384
16	BUILD_YEAR	0.002287
17	BUILD_AGE	0.002267
15	SALE_MONTH	0.002081
5	N_ROOM	0.001705

### 3. APPRENDIMENTO SUPERVISIONATO

L'apprendimento supervisionato è una tecnica di machine learning dove un modello impara a fare previsioni a partire da un dataset etichettato. In questo contesto, ogni esempio di training è composto da un input (le caratteristiche o feature) e un output noto (il target), che il modello deve imparare a predire. Durante l'addestramento, l'algoritmo analizza queste coppie input-output e cerca di identificare un pattern o una funzione che possa collegarli.

L'obiettivo è addestrare il modello in modo che, una volta addestrato, sia in grado di predire correttamente il target anche per dati mai visti prima.

L'apprendimento supervisionato è utilizzato in problemi di *classificazione*, e di *regressione*, e la scelta va fatta in base al tipo di target a disposizione, **nel nostro caso il target è continuo**.

Le tipologie di target, in base alla loro natura si affrontano con diverse tecniche di apprendimento:

- **Target Continuo:** un target è continuo quando può assumere un *numero infinito di valori in un intervallo specifico*. Esempi comuni sono valori come il prezzo di una casa o la temperatura, che possono variare in modo fluido e non si limitano a una serie di categorie.

I problemi con target continuo vengono trattati tramite *modelli di regressione*, dove l'obiettivo è predire un valore numerico preciso.

- **Target Discreto:** un target è discreto quando può assumere un *numero finito di categorie o classi*. Per esempio, nel caso della classificazione di email come "spam" o "non spam", i possibili valori target sono solo due categorie. Altri esempi sono la classificazione del colore di un'auto (ad esempio, rosso, blu, verde) o l'esito di un test medico (positivo o negativo).

Questo tipo di problema si risolve tramite *modelli di classificazione* che predicono categorie ben definite.

- **Target Strutturato:** un target è strutturato quando il valore di output è rappresentato da una *struttura complessa, come un vettore, una sequenza o un grafo*. Ad esempio, nel caso della traduzione automatica, il target è una frase in un'altra lingua, che è una sequenza di parole e non un singolo valore. Un altro esempio è il riconoscimento di oggetti nelle immagini, dove il modello deve restituire le coordinate dei confini di ciascun oggetto (tipicamente un problema di output strutturato).

I problemi con target strutturato si risolvono tramite tecniche speciali, come le *reti neurali ricorrenti (RNN)* per le sequenze o le *reti neurali convoluzionali (CNN)* per immagini, oppure con modelli come *Conditional Random Fields (CRF)* e *Hidden Markov Models (HMM)*.

Le fasi vengono suddivise in:

- scelta degli iper-parametri
- fase di addestramento
- fase di test
- valutazione delle prestazioni

### 3.1 Scelta dei modelli

L'obiettivo quindi è predire il prezzo della casa basandosi sulle caratteristiche influenti e come l'area, l'anno di costruzione, numero di stanze ecc. I modelli scelti sono DecisionTree, RandomForest, Lasso, Ridge e GradientBoosting, ciascuno adatto alla regressione.

- Il **Decision Tree Regressor** costruisce un albero basato su condizioni di soglia applicate ai valori delle feature. Partendo dalla radice, ogni nodo rappresenta una condizione sui dati (condizioni if), quindi può essere true o false, e la risposta determina il percorso da seguire nell'albero. Ogni decisione suddivide il dataset in gruppi sempre più piccoli e omogenei



rispetto alla variabile target, fino ad arrivare alle foglie, che contengono le previsioni finali (la risposta). In questo modo, un decision tree modella il processo decisionale attraverso una sequenza di scelte, che formano un cammino dalle radici alle foglie per ciascuna richiesta data di input.

È preferibile un albero decisionale più piccolo, cioè con meno nodi o meno profondo, che sia comunque coerente con il training set. Dato che esplorare tutto lo spazio possibile degli alberi di decisione è impraticabile, si utilizzano tecniche greedy per costruire l'albero minimizzando una determinata loss

- Il **Random Forest Regressor**, utilizza molti decision tree addestrati su campioni casuali del dataset, migliorando la robustezza del modello. Questa tecnica, chiamata bagging, introduce diversità nei dati di addestramento. A ogni nodo, il modello sceglie casualmente un sottoinsieme di feature, riducendo la correlazione tra gli alberi e migliorando la generalizzazione sui dati di test. Le previsioni finali del modello sono ottenute facendo la media delle previsioni di tutti gli alberi.
- Il **Gradient Boosting Regressor** costruisce una sequenza di modelli semplici, come decision tree poco profondi. Ogni modello successivo viene addestrato per correggere gli errori del precedente. Si calcolano gli errori e un nuovo modello viene addestrato per predire questi residui. La previsione complessiva viene aggiornata sommando le previsioni di tutti i modelli. Questo processo iterativo continua fino a ridurre al minimo l'errore complessivo della previsione. Così facendo, il Gradient Boosting migliora progressivamente la precisione del modello.
- La **Lasso Regression** (Least Absolute Shrinkage and Selection Operator) è una forma di regressione lineare che include una penalizzazione L1 sui coefficienti del modello, la quale tende a ridurre a zero i coefficienti meno rilevanti. Questa penalizzazione L1 ha un effetto di sparse selection (selezione sparsa), che elimina in modo efficace le feature che non contribuiscono significativamente alla previsione. Grazie a questa caratteristica, Lasso è particolarmente utile quando si lavora con dataset che contengono molte feature irrilevanti o ridondanti, rendendo il modello più interpretabile e meno soggetto a sovradattamento. Il modello finale di Lasso risulta quindi più semplice, con solo le feature essenziali, ma preserva una buona capacità predittiva.
- La **Ridge Regression** è un'altra variante della regressione lineare che applica una penalizzazione L2 ai coefficienti del modello, riducendone i valori estremi senza però eliminarli. Questa penalizzazione L2 è utile per evitare che i coefficienti diventino troppo grandi, riducendo la varianza del modello e migliorando la stabilità in presenza di feature correlate. La Ridge Regression è efficace in situazioni in cui tutte le feature hanno un ruolo, anche minimo, nella previsione, poiché tende a preservare tutte le variabili riducendo i pesi dei coefficienti. Questa stabilità rende Ridge adatta per dataset in cui le feature sono strettamente correlate, impedendo al modello di adattarsi eccessivamente a una singola feature e bilanciando l'importanza delle variabili.

### 3.2 Scelta degli iper-parametri

L'importanza e l'efficacia delle prestazioni dei modelli variano in base ai parametri impostati per i modelli. Per condurre questo task viene utilizzata la tecnica GridSearch, inclusa nella libreria Scikit-learn, che si occupa di trovare la combinazione ottimale di iperparametri per un modello di machine learning. Funziona esplorando sistematicamente un insieme predefinito di iperparametri, addestrando e valutando il modello per ogni combinazione possibile.

Per la scelta degli iper-parametri di ogni modello si utilizza la Grid Search **con 5-Fold Cross Validation ripetuta 2 volte**, in modo da esplorare lo spazio degli iper-parametri assicurandosi che i risultati siano effettivamente validi e non frutto di casualità. È stato scelto questo numero di fold per il grande numero di esempi nel dataset con l'obiettivo di mantenere la durata dell'esecuzione entro un periodo di tempo ragionevole.

La metrica utilizzata per i modelli è *neg\_mean\_squared\_error*

```
repeated_kfold = RepeatedKFold(n_splits=5, n_repeats=2, random_state=42)
```

```
for model_name, model in models.items():
    print(f"Addestramento {model_name} con RepeatedKFold...")
    start_train_time = time.time()
    grid_search = GridSearchCV(
        model,
        param_grids[model_name],
        cv=repeated_kfold, # Usa RepeatedKFold come metodo di validazione
        scoring='neg_mean_squared_error'
    )
    grid_search.fit(X_train, y_train)
```

*grid\_search.fit(X\_train, y\_train)* addestra il modello sui dati di addestramento, esplorando tutte le combinazioni degli iperparametri specificati in *param\_grids[model\_name]*

#### Iperparametri utilizzati per ogni modello:

Decision tree regression:

```
'DecisionTree': {
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10]
},
```

- **max\_depth**: Limita la profondità dell'albero, evitando che cresca eccessivamente e si adatti troppo ai dati di training, riducendo così l'overfitting.
- **min\_samples\_split**: Controlla la suddivisione dei nodi, con valori più alti che impediscono la creazione di nodi molto specifici e quindi riducono il rischio di overfitting.

Randomforest regression

```
param_grids = {
    'RandomForest': {
        'n_estimators': [100, 200],
        'max_depth': [None, 10, 20],
        'min_samples_split': [2, 5, 10]
    },
    ...
}
```

- **n\_estimators:** Questo parametro indica il numero di decision tree nel Random Forest. Più alberi in generale portano a una maggiore stabilità del modello, ma aumentano anche il tempo di addestramento.
- **max\_depth:** Limita la profondità massima degli alberi, il che riduce la complessità del modello e può aiutare a prevenire l'overfitting.
- **min\_samples\_split:** Specifica il numero minimo di campioni richiesti per suddividere un nodo. Valori più alti riducono la complessità dell'albero, dato che impongono più campioni per ogni decisione e quindi aumentano la generalizzazione.

### GradientBoosting Regressor

```
'GradientBoosting': {
    'n_estimators': [100, 200],
    'learning_rate': [0.01, 0.1, 0.2],
    'max_depth': [3, 5, 10]
}
```

- **n\_estimators:** È il numero totale di alberi deboli (o stadi) costruiti nella sequenza. Un numero maggiore consente una rappresentazione più dettagliata degli errori, ma richiede un bilanciamento con `learning_rate` per evitare l'overfitting.
- **learning\_rate:** Determina quanto ciascun albero contribuisce al modello finale. Un valore più basso rende il modello più stabile ma richiede più alberi per raggiungere una buona accuratezza.
- **max\_depth:** Limita la profondità di ciascun albero all'interno del boosting, con alberi più bassi che riducono il rischio di sovradattamento, in particolare nei modelli di boosting, che tendono a essere sensibili all'overfitting quando gli alberi sono troppo complessi.

### Lasso

```
'Lasso': {
    'alpha': [0.01, 0.1, 1, 10]
}
```

- **alpha:** È il parametro di regolarizzazione L1, che controlla l'entità della penalizzazione applicata ai coefficienti. Valori più alti per `alpha` aumentano l'intensità della regolarizzazione, eliminando coefficienti irrilevanti e producendo un modello più semplice, ma al costo di una possibile perdita di precisione.

### Ridge

```
'Ridge': {
    'alpha': [0.01, 0.1, 1, 10]
},
```

- **alpha:** Anche qui, `alpha` è il parametro di regolarizzazione (in questo caso L2), che riduce i coefficienti e stabilizza il modello, limitando la varianza. Come in Lasso, valori più alti

aumentano la regolarizzazione e riducono l'importanza delle feature meno rilevanti senza eliminarle.

### 3.3 Valutazione dei modelli

#### **Addestramento RandomForest...**

```
Migliori iperparametri per RandomForest: {'max_depth': None,
'min_samples_split': 5, 'n_estimators': 200}
RandomForest MSE: 561440183312.3574
RandomForest MAE: 596192.1844419925
RandomForest R2: 0.9572639906290961
```

#### **Addestramento Lasso...**

```
Migliori iperparametri per Lasso: {'alpha': 10}
Lasso MSE: 1901049798438.0662
Lasso MAE: 1103898.1108480382
Lasso R2: 0.8552948570205839
```

#### **Addestramento Ridge...**

```
Migliori iperparametri per Ridge: {'alpha': 1}
Ridge MSE: 1901082187476.7986
Ridge MAE: 1103906.0907978127
Ridge R2: 0.8552923916141097
```

#### **Addestramento DecisionTree...**

```
Migliori iperparametri per DecisionTree: {'max_depth': 10,
'min_samples_split': 10}
DecisionTree MSE: 847455966686.2823
DecisionTree MAE: 707614.7326468945
DecisionTree R2: 0.9354928855999177
```

#### **Addestramento GradientBoosting...**

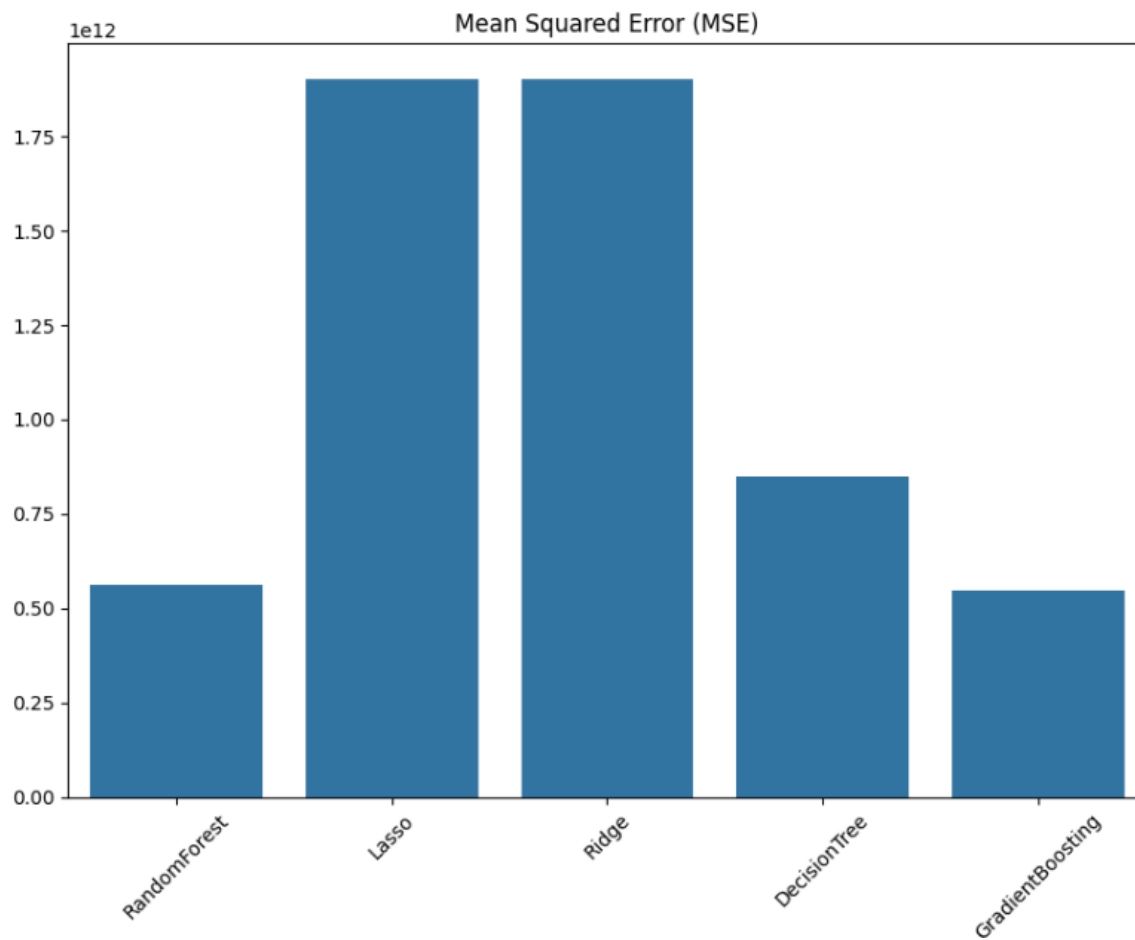
```
Migliori iperparametri per GradientBoosting: {'learning_rate': 0.1,
'max_depth': 5, 'n_estimators': 200}
GradientBoosting MSE: 545096985830.48944
GradientBoosting MAE: 579090.6317725645
GradientBoosting R2: 0.9585080110278765
```

Per valutare i modelli ho usato le metriche, a Mean Squared Error (MSE), Mean Absolute Error (MAE) e al punteggio  $R^2$ , che misurano rispettivamente l'accuratezza delle previsioni, la precisione complessiva e la bontà dell'adattamento. Dopo l'addestramento, il modello con il minor MSE viene selezionato come migliore e salvato per l'uso futuro.

**L'errore quadratico medio (MSE)** misura la media dei quadrati degli errori, ovvero la differenza tra i valori previsti dal modello e i valori reali.

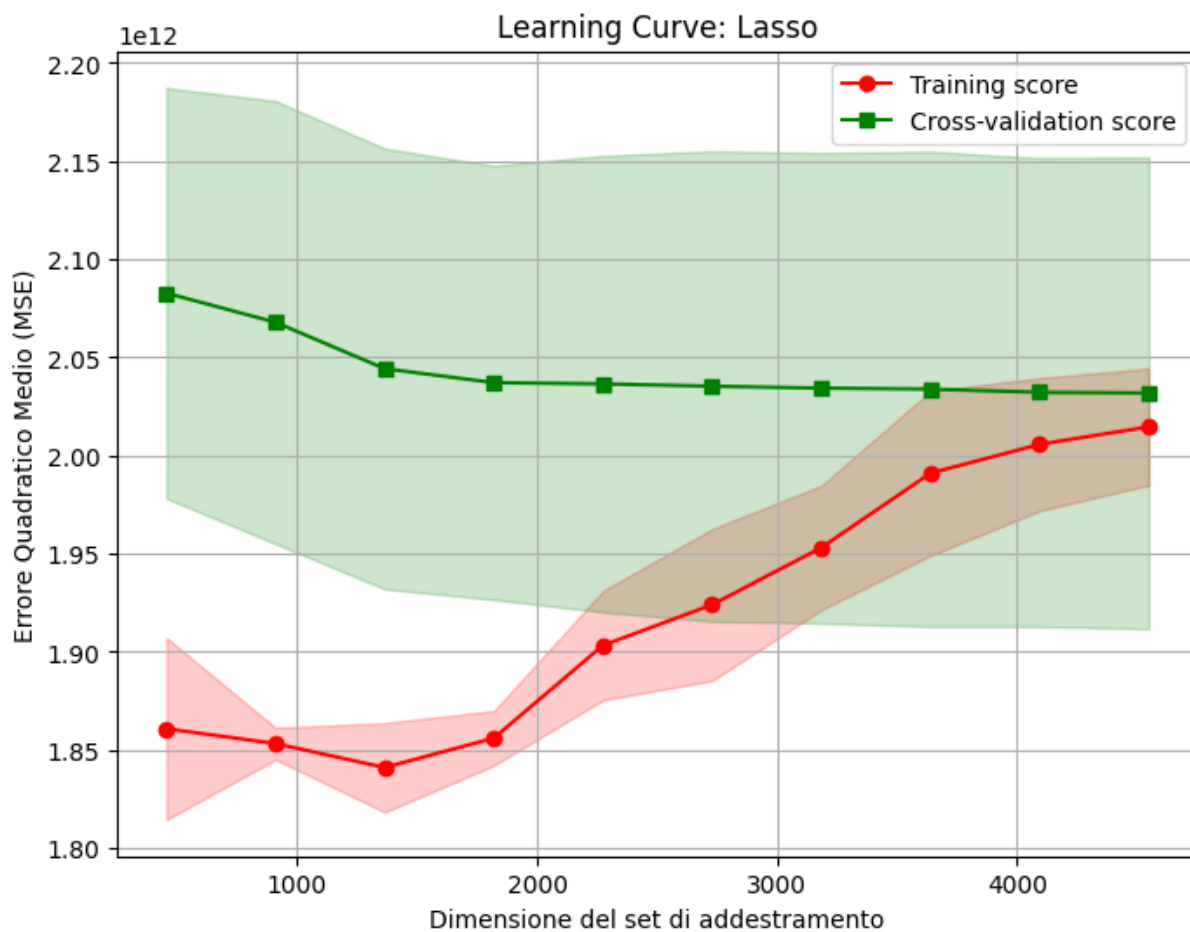
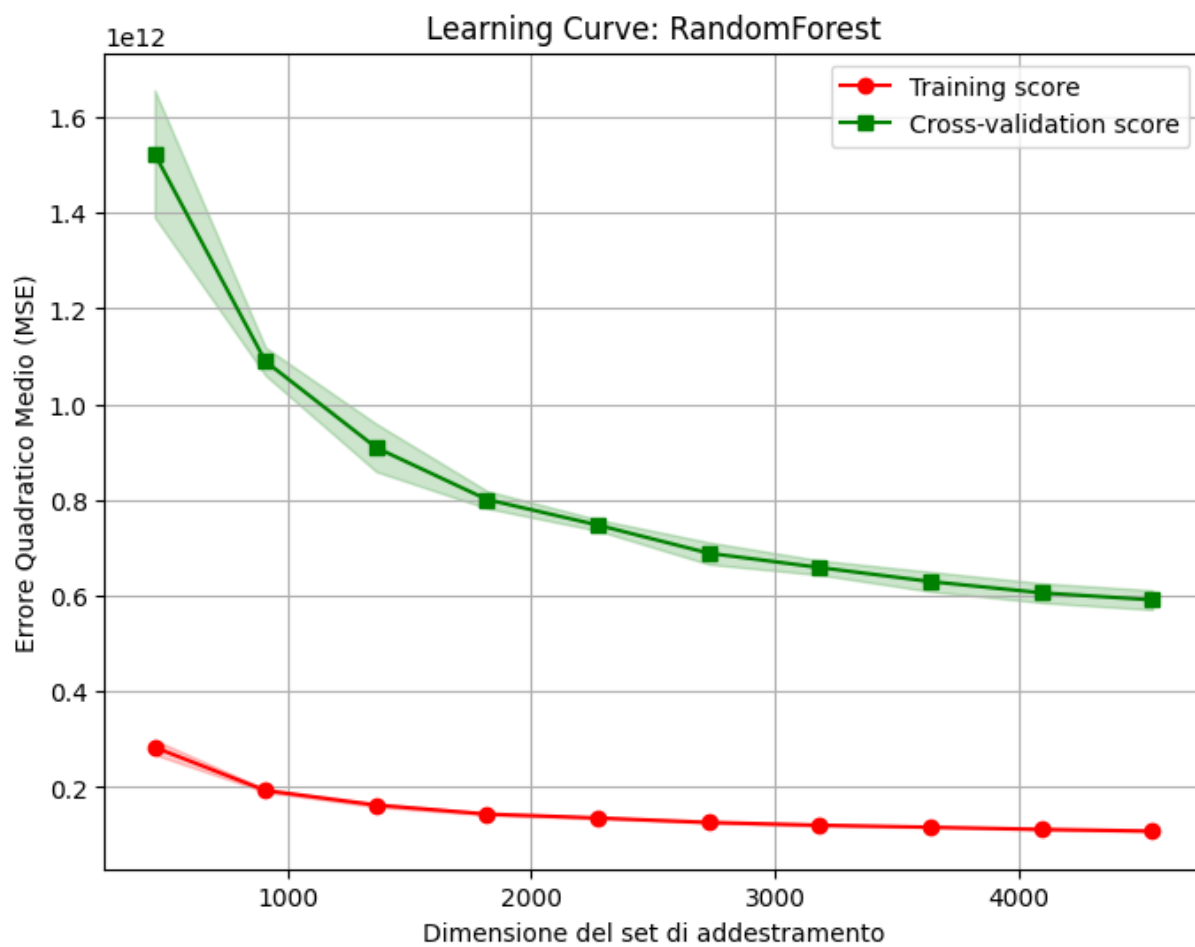
**L'errore assoluto medio (MAE)** è la media dei valori assoluti delle differenze tra i valori previsti e quelli reali.

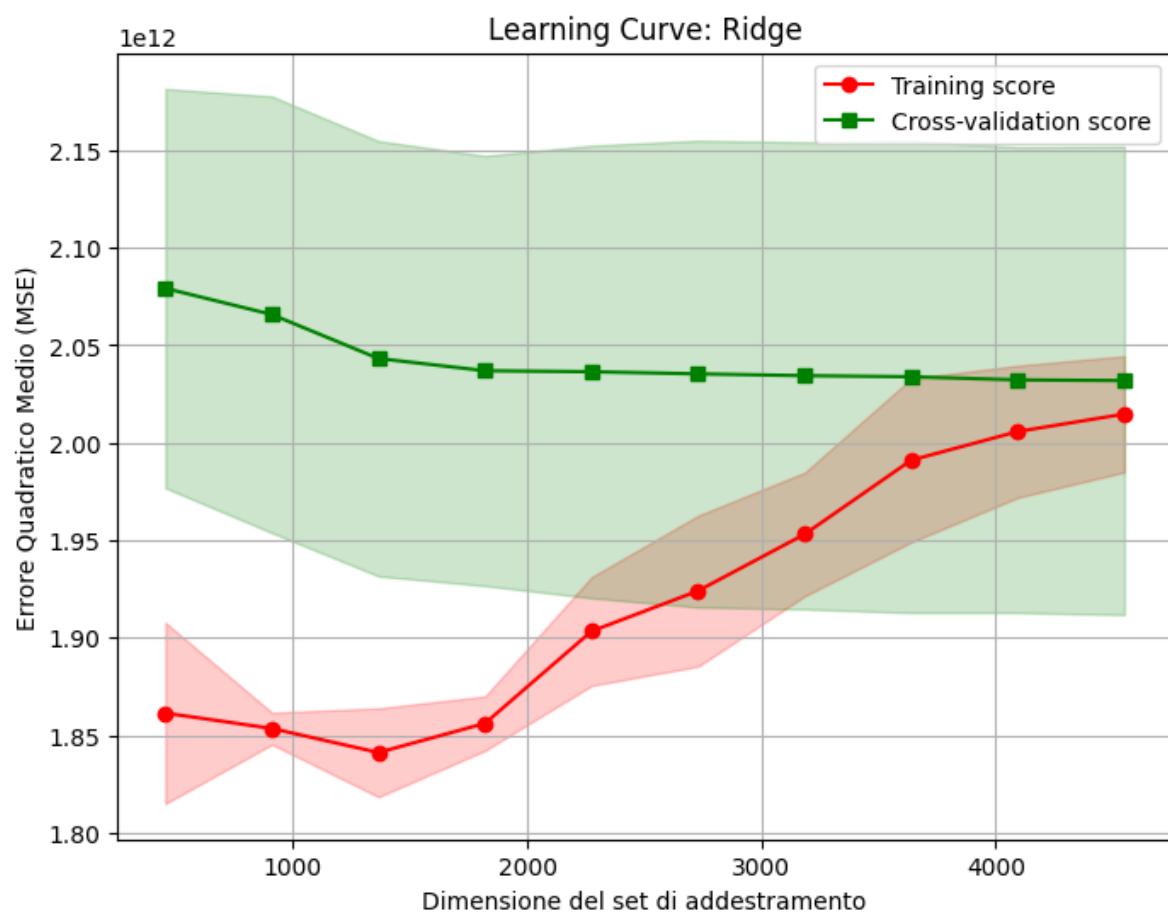
**Il coefficiente di determinazione  $R^2$**  misura la percentuale di varianza nei dati che è spiegata dal modello di regressione.

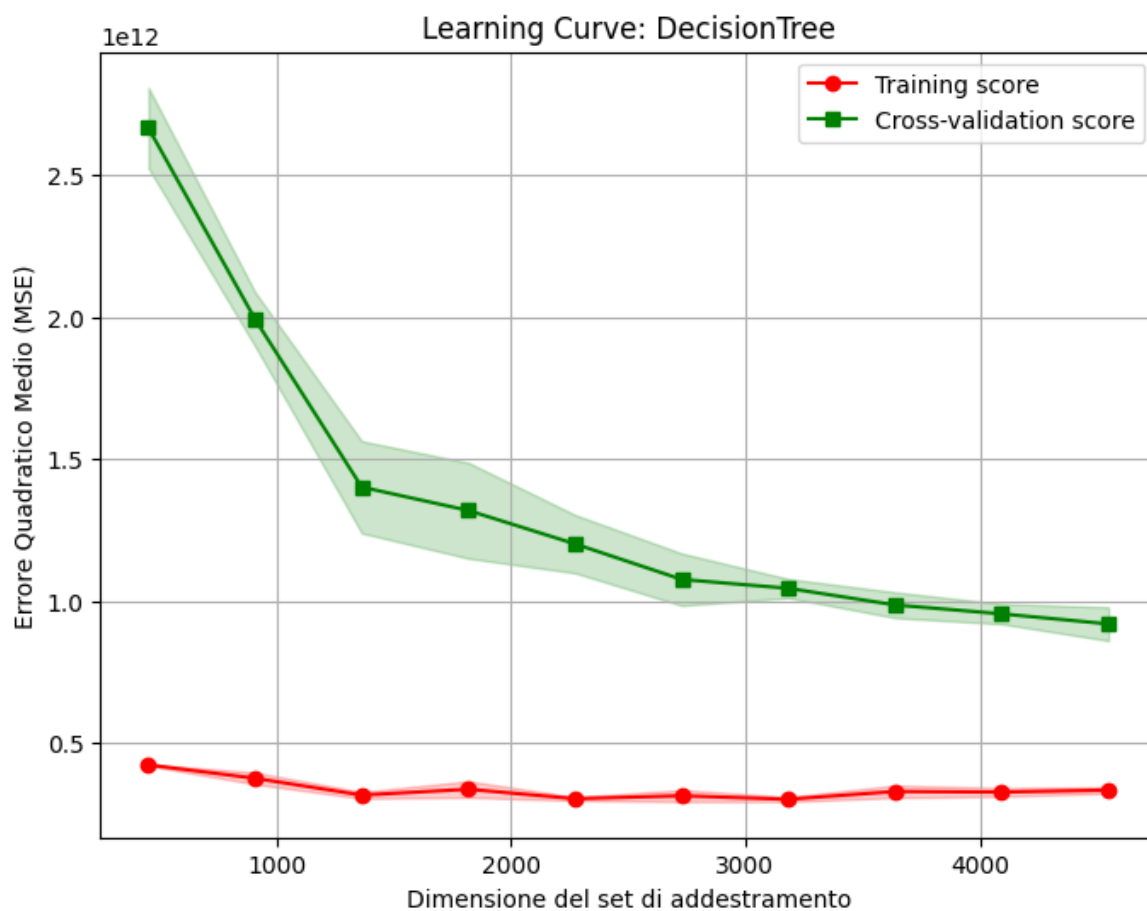


Come si può notare GradientBoosting ha l'errore MSE e  $R^2$  più basso(0.9585) tra tutti i modelli, suggerendo prestazioni complessivamente eccellenti e una maggiore accuratezza nelle predizioni rispetto agli altri modelli.

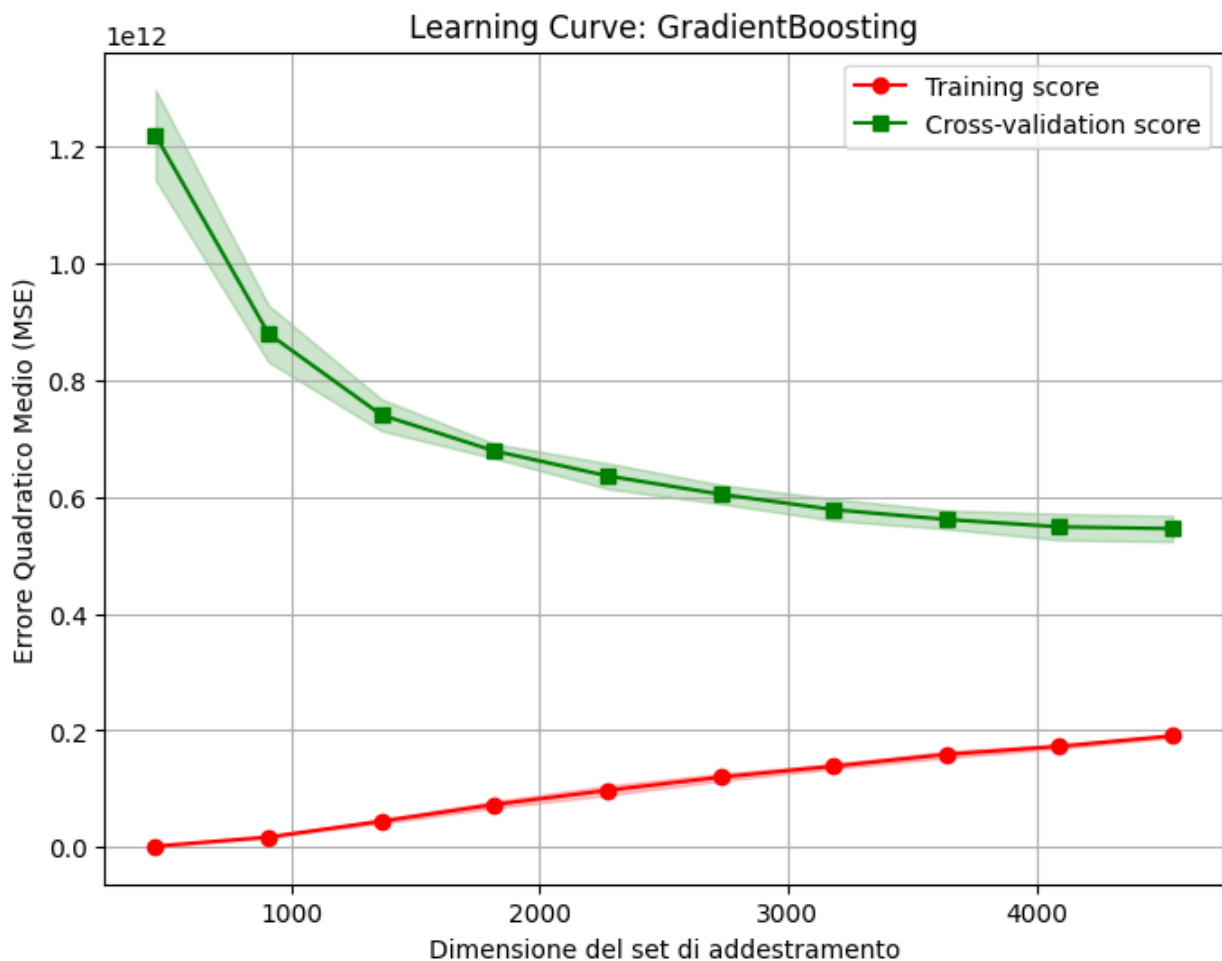
Ora andiamo ad analizzare le curve di apprendimento per ogni modello.











Dai grafici analizzati, emerge che i modelli Random Forest e Gradient Boosting generalmente presentano prestazioni migliori rispetto a Lasso, Ridge e Decision Tree. I primi due mostrano una minore tendenza all'overfitting e una maggiore capacità di generalizzare i dati.

### 3.4 Varianza e deviazioni standard

La **varianza** è una misura della dispersione dei dati rispetto alla media. Indica quanto i valori di un insieme di dati differiscano tra loro e dalla media del dataset. Matematicamente, è calcolata come la media dei quadrati delle differenze di ciascun valore dalla media del dataset. Una varianza più alta indica che i dati sono più dispersi, mentre una varianza più bassa indica che i dati sono più concentrati attorno alla media.

La **deviazione standard** è la radice quadrata della varianza e rappresenta anch'essa la dispersione dei dati rispetto alla media, ma è espressa nella stessa unità dei dati originali. È utile per capire quanto, mediamente, i valori si discostano dalla media. Poiché è calcolata come radice quadrata, risulta più intuitiva della varianza in termini pratici di dispersione dei valori.

Deviazioni standard e varianze per il modello RandomForest:

Train Size	Train STD	Train VAR	Test STD	Test VAR
454	1.78418	3.18330	1.250608	1.564020
909	3.20161	1.02503	4.521985	2.044834
1364	6.47559	4.19332	6.243432	3.898044
3184	1.89019	3.57282	1.946891	3.790386
3639	1.41691	2.00765	2.107380	4.441051
4094	2.32292	5.39596	2.000350	4.001401
4549	1.76121	3.10188	2.230332	4.974381

Deviazioni standard e varianze per il modello Lasso:

Train Size	Train STD	Train VAR	Test STD	Test VAR
454	4.63169	2.14525	1.045591	1.093261
909	8.09141	6.54710	1.127972	1.272320
1364	2.27145	5.15950	1.122022	1.258934
3184	3.16637	1.00259	1.197519	1.434053
3639	4.19354	1.75857	1.210064	1.464255
4094	3.38036	1.14268	1.193803	1.425165
4549	2.97384	0.88437	1.200750	1.441800

Deviazioni standard e varianze per il modello Ridge:

Train Size	Train STD	Train VAR	Test STD	Test VAR
454	4.63934	2.15234	1.023098	1.046729
909	8.14278	6.63049	1.118771	1.251649
1364	2.26978	5.15193	1.115515	1.244374
3184	3.16630	1.00254	1.196994	1.432795
3639	4.19350	1.75854	1.209361	1.462555
4094	3.38033	1.14266	1.193109	1.423509
4549	2.97382	0.88436	1.200284	1.440682

Deviazioni standard e varianze per il modello DecisionTree:

Train Size	Train STD	Train VAR	Test STD	Test VAR
454	3.29432	1.08525	1.628922	2.653387
909	2.06824	4.27764	1.018728	1.037807
1364	1.16994	1.36876	1.749138	3.059483
3184	0.86763	0.75279	0.415740	1.728400
3639	2.15955	4.66366	0.465747	2.169204
4094	1.49587	2.23765	0.445034	1.980553
4549	1.12669	1.26943	0.625878	3.917242

Deviazioni standard e varianze per il modello GradientBoosting:

Train Size	Train STD	Train VAR	Test STD	Test VAR
454	0.10939	0.01196	0.825630	6.816661
909	0.72199	0.52126	0.481629	2.319666
1364	3.41405	1.16558	0.266280	0.709054
3184	4.41724	1.95120	0.186704	0.348583
3639	5.78721	3.34919	0.146404	0.214342
4094	3.49325	1.22028	0.226014	0.510824
4549	2.75178	0.75722	0.224881	0.505718

La **varianza** del training in tutti i modelli tende generalmente a essere inferiore a quella del test, il che è comune e indica che i modelli si adattano meglio ai dati di addestramento.

La varianza elevata nel test rispetto al training suggerisce una difficoltà del modello a generalizzare, evidenziando la possibilità di overfitting.

In generale, *GradientBoosting* appare più stabile e meno incline a varianza elevata, rendendolo potenzialmente migliore per questo problema.

Soluzioni per gli altri modelli: *DecisionTree* e *RandomForest* mostrano un'alta varianza, in particolare su set più piccoli, e potrebbero beneficiare di tecniche di pruning o del tuning di parametri per migliorare la generalizzazione, invece *Lasso* e *Ridge* presentano una buona costanza, ma potrebbero migliorare con regularizzazione per ridurre ulteriormente la varianza nel test.

#### Calcolo del prezzo con il miglior modello (gradientboosting)

Vengono richiesti in input all'utente alcuni dati delle feature, successivamente viene creato un dataframe e viene passato al modello i parametri inseriti e tramite la funzione `predici_prezzo`, il modello calcolerà il valore di output(prezzo). Esempio di output generato

```
- Area: Adyar
- Tipo di Edificio: House
- Facilità di Parco: No
- Metri Quadrati Interni: 3376.92
- Numero di Camere da Letto: 2
- Numero di Bagni: 3
- Totale Stanze: 6
- Qualità Stanze (0-10): 2.28
- Qualità Bagno (0-10): 3.32
- Qualità Camera (0-10): 0.59
- Qualità Complessiva (0-10): 3.75
- Tassa di Registrazione (INR): 12710.13
- Commissione (INR): 5119.25
- Anno di Vendita: 2010
- Anno di Costruzione: 1989
- Mese di Vendita: 12
```

```
### Prezzo Previsto della Casa: 61,956.56 Euro
```

## 4. BASE DI CONOSCENZA E PROLOG

### 4.1 Conversione da albero a prolog

Per il ragionamento logico viene costruita una cosiddetta Knowledge Base, contenete gli assiomi(ovvero affermazioni sempre vere),costituiti da fatti e regole. I fatti indicano verità immutabili, invece una regola è una dichiarazione che afferma che qualcosa è vero in base ad altre cose che sono vere.

Successivamente l'albero 'decision tree' precedentemente addestrato, viene convertito in prolog, linguaggio di programmazione dichiarativo basato sul ragionamento logico, per creare fatti e regole. Libreria usata `pyswip` .

### 4.2 Regole e fatti

Nel progetto sono state create 3 regole e fatti 2.

- Un **fatto** di tipo 'leaf(NodeID, Prediction)' viene generato per ogni nodo foglia dell'albero di decisione.

Questo rappresenta un nodo terminale dell'albero dove viene effettuata una previsione (il valore della foglia).

- Un **fatto** di tipo 'node(NodeID, FeatureIndex, Threshold, LeftChild, RightChild)' viene creato per ogni nodo interno dell'albero.

Questo rappresenta un nodo di decisione, ossia un punto in cui il percorso di previsione è suddiviso in base al valore di una caratteristica.

- **Regola** principale 'predire\_prezzo(Features, Prezzo)' che utilizza una chiamata ricorsiva 'percorri\_albero' per attraversare l'albero di decisione fino a raggiungere una foglia e ottenere unaa previsione.

Questa regola:

- Fa un confronto tra i valori di soglia e il valore di input delle caratteristiche.
- Decide in quale direzione muoversi (figlio sinistro o destro) per arrivare alla previsione.
- `predire_prezzo(Features, Prezzo) :-  
    percorri_albero(0, Features, Prezzo).`
- `percorri_albero(NodeID, _, Predizione) :-  
    leaf(NodeID, Predizione).`
- `percorri_albero(NodeID, Features, Predizione) :-  
    node(NodeID, FeatureIndex, Threshold, LeftChild, RightChild),  
    nth0(FeatureIndex, Features, FeatureValue),  
    (FeatureValue =< Threshold ->  
        percorri_albero(LeftChild, Features, Predizione)  
        percorri_albero(RightChild, Features, Predizione)).`

## ESEMPIO DI OUTPUT

Predizione Prolog: 2156875.0

Predizione Decision Tree: 2156875.0

Features usate per Prolog: [5, 4337, 1, 2, 4, 2, 1, 9, 7, 3, 3, 28833.12845182201, 3050.8218212634033, 2004, 2005, 2, 19, 0]

## 5. APPRENDIMENTO BAYESIANO

Per questo argomento sono state utilizzate la libreria Python pgmpy per la gestione della rete bayesiana e le librerie matplotlib e networkx per la visualizzazione grafica.

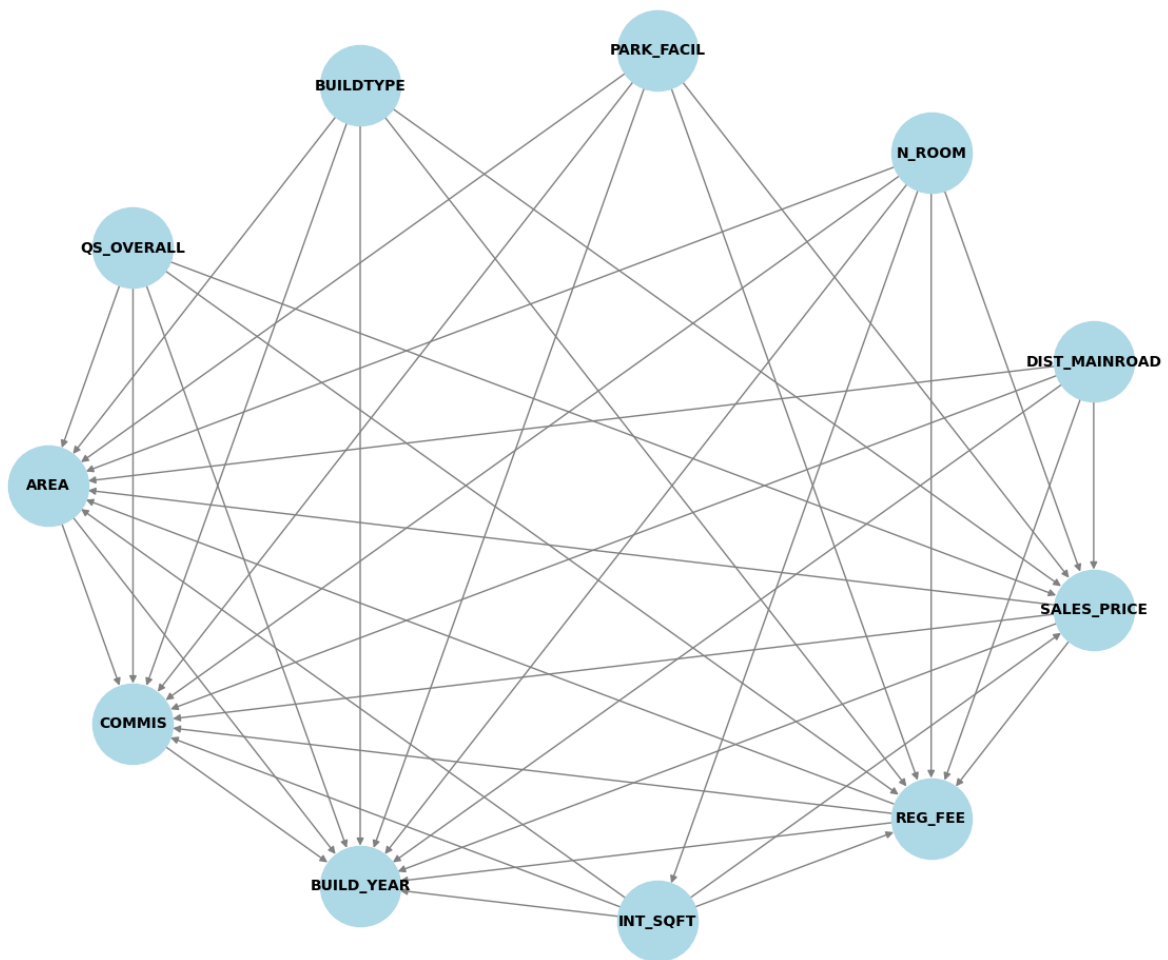
L'apprendimento bayesiano consente di creare una rete bayesiana, che può essere utilizzata come classificatore per predire il prezzo di vendita o per effettuare previsioni anche in assenza di valori per alcune variabili.

### 5.1. Struttura della rete bayesiana

Per costruire una rete bayesiana, è necessario definirne la struttura, ossia le variabili e le relazioni di dipendenza tra esse. La struttura può essere definita manualmente oppure appresa dai dati, come in questo caso.

L'algoritmo HillClimbSearch è stato utilizzato per apprendere la struttura. Questo metodo di ricerca locale modifica iterativamente il grafo fino a trovare una configurazione che rappresenta in modo appropriato le dipendenze tra le variabili.

Struttura della rete bayesiana appresa



La rete bayesiana risultante evidenzia le connessioni dirette tra le variabili principali, come il collegamento tra caratteristiche come AREA e INT\_SQFT e la variabile target SALES\_PRICE.

## 5.2. Addestramento della rete bayesiana

La rete bayesiana è stata addestrata basandosi sulla struttura definita in precedenza, utilizzando le Conditional Probability Distributions (CPD). Ogni CPD definisce la distribuzione di probabilità di una variabile in funzione delle variabili correlate.

CPD per variabile SALES\_PRICE

BUILDTYPE	BUILDTYPE(0.0)	...	BUILDTYPE(4.0)	
DIST_MAINROAD	DIST_MAINROAD(0.0)	...	DIST_MAINROAD(4.0)	
INT_SQFT	INT_SQFT(0.0)	...	INT_SQFT(4.0)	
N_ROOM	N_ROOM(0.0)	...	N_ROOM(4.0)	
PARK_FACIL	PARK_FACIL(0.0)	...	PARK_FACIL(4.0)	
QS_OVERALL	QS_OVERALL(0.0)	...	QS_OVERALL(4.0)	
SALES_PRICE(0.0)	0.003720927798089027	...	0.19999999999999998	
SALES_PRICE(1.0)	0.003720927798089027	...	0.19999999999999998	
SALES_PRICE(2.0)	0.9851162888076438	...	0.19999999999999998	
SALES_PRICE(3.0)	0.003720927798089027	...	0.19999999999999998	
SALES_PRICE(4.0)	0.003720927798089027	...	0.19999999999999998	

Si può osservare come una casa con una superficie interna maggiore (**INT\_SQFT**) e un numero elevato di stanze (**N\_ROOM**) abbia una maggiore probabilità di avere un prezzo di vendita (**SALES\_PRICE**) più alto, mentre caratteristiche come una distanza maggiore dalla strada principale (**DIST\_MAINROAD**) tendono a ridurre la probabilità di un prezzo elevato.

### 5.3 Utilizzo della rete bayesiana

La rete bayesiana può essere utilizzata sia come un classificatore per prevedere la variabile target **SALES\_PRICE** sulla base di altre variabili osservabili, sia per inferire valori mancanti di altre variabili.

```
Evidence: {'AREA': 2.0, 'INT_SQFT': 1.0, 'DIST_MAINROAD': 2.0, 'N_ROOM': 1.0, 'PARK_FACIL': 4.0, 'BUILDTYPE': 4.0, 'QS_OVERALL': 2.0, 'REG_FEE': 0.0, 'COMMIS': 0.0, 'BUILD_YEAR': 2.0}
```

```
Predicted variables: {'SALES_PRICE': 1.0}
```

```
Evidence: {'AREA': 3, 'INT_SQFT': 2}
```

```
Predicted variables: {'SALES_PRICE': 2.0}
```

```
Evidence from row: {'AREA': 1.0, 'INT_SQFT': 1.0, 'DIST_MAINROAD': 1.0, 'N_ROOM': 2.0, 'PARK_FACIL': 4.0, 'BUILDTYPE': 0.0, 'QS_OVERALL': 3.0, 'REG_FEE': 1.0, 'COMMIS': 1.0, 'SALES_PRICE': 2.0, 'BUILD_YEAR': 2.0}
```

## **6. CONCLUSIONI**

Complessivamente in questo progetto sono stati sviluppati modelli di apprendimento supervisionato per prevedere il prezzo di vendita di case/appartamenti basandoci su dataset online,

Il modello di Gradient Boosting ha mostrato le migliori performance in termini di errore medio assoluto (MAE), seguito da Random Forest e Decision Tree.

I modelli di regressione lineare hanno mostrato prestazioni inferiori rispetto agli altri.

È stata creata una base di conoscenza per stimare il prezzo di una casa, usando un albero decisionale o un calcolo “manuale” basandosi dai valori delle caratteristiche.

## **POSSIBILI SVILUPPI**

Per sviluppare ulteriormente il progetto, si potrebbero testare modelli di apprendimento supervisionato più sofisticati per migliorare ulteriormente l'accuratezza delle predizioni.

Si potrebbe introdurre la grafica per poter prendere i dati di input dagli utenti per utilizzare le predizioni.

Si possono creare nuove applicazioni basate su intelligenza artificiale per poter sfruttare a pieno le ontologie e web semantico, con il ragionamento automatico.

## **RIFERIMENTI BIBLIOGRAFICI**

Apprendimento supervisionato: D. Poole, A. Mackworth: Artificial Intelligence: Foundations of Computational Agents. 3rd ed. Cambridge University Press [Ch.7]

Basi di conoscenza: D. Poole, A. Mackworth: Artificial Intelligence: Foundations of Computational Agents. 3rd ed. Cambridge University Press [Ch.15]