

# Setup Iniziale e Librerie

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # Sklearn Preprocessing & Split
7 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
8 from sklearn.preprocessing import StandardScaler, OneHotEncoder, OrdinalEncoder
9 from sklearn.impute import SimpleImputer
10 from sklearn.compose import ColumnTransformer
11 from sklearn.pipeline import Pipeline
12
13 # Sklearn Models
14 from sklearn.linear_model import LinearRegression, LogisticRegression
15 from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
16 from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
17 from sklearn.neighbors import KNeighborsClassifier
18 from sklearn.cluster import KMeans, DBSCAN, AgglomerativeClustering
19
20 # Sklearn Metrics
21 from sklearn.metrics import mean_squared_error, r2_score, silhouette_score
22 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix,
    ConfusionMatrixDisplay
```

Listing 1: Tutti gli Import necessari

Il Preprocessing è la fase in cui trasformiamo i dati grezzi in un formato comprensibile per gli algoritmi di Machine Learning (che "ragionano" solo con i numeri).

## 1. Gestione dei Valori Mancanti (Missing Values)

**Teoria:** Gli algoritmi di 'sklearn' non accettano NaN (Not a Number). Dobbiamo decidere se eliminare le righe (che però riduce il dataset) o riempirle (imputazione).

**Requisito d'Esame:** Spesso viene richiesto di imputare la media per i numeri e una stringa fissa per le categorie.

```
1 from sklearn.impute import SimpleImputer
2
3 # Per le colonne numeriche: usiamo la media
4 # Strategy='mean' riempie i buchi con la media della colonna
5 num_imputer = SimpleImputer(strategy='mean')
6
7 # Per le colonne categoriche: usiamo un valore fisso
8 # Strategy='constant' mette un valore standard dove manca il dato
9 cat_imputer = SimpleImputer(strategy='constant', fill_value='unknown')
```

Listing 2: Imputazione dei dati

## Approfondimento: Pandas vs Sklearn e il Data Leakage

Esiste un modo più veloce per riempire i valori mancanti usando direttamente Pandas, ma comporta dei rischi metodologici.

```
1 # Metodo veloce: calcola la media su TUTTA la colonna e riempie i buchi
2 # Rischio: Se fatto prima dello split Train/Test, introduci Data Leakage!
3 df['colonna'] = df['colonna'].fillna(df['colonna'].mean())
```

Listing 3: Metodo rapido con Pandas (Attenzione al Leakage!)

## Cos'è il Data Leakage?

Il **Data Leakage** (fuga di dati) avviene quando informazioni del *Test Set* (che il modello non dovrebbe mai vedere) influenzano il training.

Se calcoliamo la media su tutto il dataset prima di dividere:

1. La media usata per riempire i buchi nel Train Set contiene "tracce" dei valori del Test Set.

2. Il modello avrà performance ottimistiche in fase di training, ma peggiori su dati reali nuovi.

### Procedura Corretta per l'Esame:

1. Dividere in `X_train` e `X_test`.
2. `imputer.fit(X_train)`: Impara la media solo dal Train.
3. `imputer.transform(X_train)` e `imputer.transform(X_test)`: Usa quella media per riempire entrambi.

## 2. Encoding (Da Testo a Numeri)

**Teoria:** I modelli matematici non sanno cosa sia "Rosso" o "Blu". Dobbiamo tradurli.

- **Ordinal Encoding:** Se c'è un ordine (es. Low < Medium < High) → 1, 2, 3.
- **One-Hot Encoding:** Se NON c'è ordine (Nominale, es. Colori, Città). Crea una colonna per ogni valore (es. `is_Red`, `is_Blue`) con 0 o 1.

```
1 from sklearn.preprocessing import OneHotEncoder
2
3 # handle_unknown='ignore': se nel test set appare una categoria mai vista, la ignora (tutti 0)
  invece di dare errore.
4 # sparse_output=False: restituisce un array numpy leggibile invece di una matrice sparsa compressa.
5 ohe = OneHotEncoder(handle_unknown='ignore', sparse_output=False)
```

Listing 4: OneHotEncoder

## 3. Feature Scaling (Mettere tutto sulla stessa scala)

**Teoria:** Molti algoritmi calcolano la "distanza" tra i punti (es. KMeans, KNN, SVM). Se una feature è in *metri* (1-10) e l'altra in *millimetri* (1000-10000), quella in millimetri dominerà il calcolo solo perché ha numeri più grandi.

Lo scaling porta tutto in un range comparabile (solitamente media 0 e varianza 1).

```
1 from sklearn.preprocessing import StandardScaler
2
3 scaler = StandardScaler()
4 # NOTA: Si fa fit() solo sul Train Set per evitare il "Data Leakage"!
5 # X_train_scaled = scaler.fit_transform(X_train)
6 # X_test_scaled = scaler.transform(X_test)
```

Listing 5: StandardScaler

## Data Visualization & Exploration (EDA)

Questa fase serve a capire la struttura dei dati prima di applicare i modelli. Le domande d'esame tipiche richiedono di individuare: *outliers*, *imbalanced distributions* e correlazioni con il target.

### 1. Boxplot: Distribuzioni e Outliers

**Teoria:** Il Boxplot riassume la distribuzione di una variabile numerica.

- **Box (Scatola):** Contiene il 50% centrale dei dati (dal 25° al 75° percentile).
- **Linea nel mezzo:** La Mediana (il valore centrale).
- **Baffi (Whiskers):** Si estendono fino a 1.5 volte l'IQR (intervallo interquartile).
- **Punti esterni:** Sono gli **Outliers** (valori anomali molto alti o bassi).

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Imposta lo stile
5 sns.set(style="whitegrid")
```

```

6
7 # Boxplot per visualizzare gli outliers
8 plt.figure(figsize=(12, 6))
9 # Orient='h' li fa orizzontali (piu' leggibili se hai tante features)
10 sns.boxplot(data=df, orient='h')
11 plt.title("Distribution & Outliers Analysis")
12 plt.show()

```

Listing 6: Generare Boxplot per tutte le feature

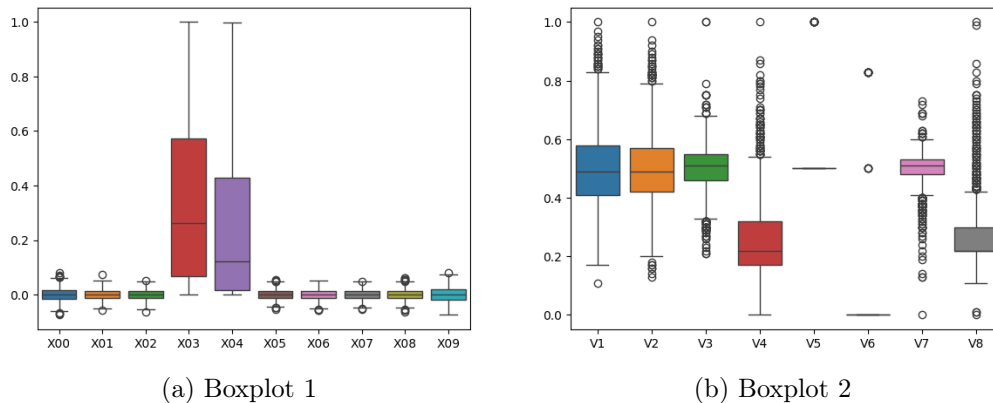


Figura 1: Confronto tra i due Boxplot, uno con molti outliers e uno con pochi

**Come commentare all'esame (Template):** *"The boxplots show that features 'Age' and 'Income' have several outliers beyond the upper whisker. The variable 'Count' shows a skewed distribution (the median is not in the center of the box)."*

## 2. Heatmap: Correlazioni

**Teoria:** Misura quanto due variabili si muovono insieme (Indice di Pearson, da -1 a 1).

- **1:** Correlazione positiva perfetta (se sale A, sale B).
- **-1:** Correlazione negativa perfetta (se sale A, scende B).
- **0:** Nessuna relazione lineare.

### Approfondimento Teorico: Indice di Pearson ( $r$ )

L'indice di Pearson misura la forza e la direzione della relazione **lineare** tra due variabili continue.

**Formula:**

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Sostanzialmente è il rapporto tra la Covarianza di X e Y e il prodotto delle loro Deviazioni Standard.

**Interpretazione dei Valori:**

- **+1 (Correlazione Positiva Perfetta):** Tutti i punti giacciono su una retta che sale. Se X aumenta, Y aumenta in proporzione fissa.
- **-1 (Correlazione Negativa Perfetta):** Tutti i punti giacciono su una retta che scende. Se X aumenta, Y diminuisce.
- **0 (Assenza di Correlazione Lineare):** Non c'è alcuna relazione *lineare*. Le variabili potrebbero essere indipendenti o legate in modo complesso.

**Perché usiamo `abs()` nel codice?** Quando cerchiamo feature predittive, ci interessa la *forza* della relazione, non la direzione. Una feature con correlazione  $-0.9$  è preziosissima (molto predittiva), tanto quanto una con  $+0.9$ . Una con  $0.01$  è inutile (rumore).

## Trappola d'Esame: Relazioni Non Lineari

**Attenzione:** Una correlazione di Pearson pari a 0 **non** significa che le variabili non sono collegate. Significa solo che non c'è una linea retta che le unisce.

*Esempio:* Se  $Y = X^2$  (una parabola perfetta), l'indice di Pearson sarà  $\approx 0$ .

**Conclusione:** Guardare sempre il grafico (scatterplot) prima di fidarsi ciecamente del numero!

**Requisito d'Esame:** Spesso si chiede di scartare le feature con correlazione assoluta col target  $< 0.15$ .

```
1 # Calcola la matrice di correlazione
2 corr_matrix = df.corr()
3
4 # Visualizza la Heatmap
5 plt.figure(figsize=(10, 8))
6 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f")
7 plt.title("Correlation Matrix")
8 plt.show()
9
10 # Esempio: Trovare feature poco correlate con il Target (es. 'price')
11 target = 'price'
12 # abs() serve per considerare sia correlazioni positive che negative
13 low_corr_features = corr_matrix.index[abs(corr_matrix[target]) < 0.15].tolist()
14 print("Features to drop:", low_corr_features)
```

Listing 7: Heatmap e Filtro Correlazione



Figura 2: più il colore è chiaro, più la correlazione è alta

**Come commentare all'esame:** *"From the heatmap, we observe that 'Feature A' has a strong positive correlation (0.85) with the target, while 'Feature B' is weakly correlated ( $< 0.15$ ) and should be dropped to reduce noise."*

## Train-Test Split (Attenzione a Stratify!)

**Teoria:** Dividiamo i dati per simulare come il modello si comporterà nel mondo reale.

- **Regressione:** Split casuale semplice.
- **Classificazione:** Usare `stratify=y` è fondamentale se le classi sono sbilanciate (es. 90% 'No', 10% 'Sì'). Garantisce che Train e Test abbiano le stesse proporzioni.

```
1 # 1. Separare Features (X) e Target (y)
2 target_name = 'price' # 0 'class', cambia in base al dataset
3 X = df.drop(columns=[target_name])
4 y = df[target_name]
5
6 # 2. Split
7 # test_size=0.2 (20% test, 80% train) lo standard
8 # random_state=42 rende tutto riproducibile
9 # stratify=y SOLO per CLASSIFICAZIONE (rimuovere per regressione)
10
11 X_train, X_test, y_train, y_test = train_test_split(
12     X, y,
13     test_size=0.2,
14     random_state=42,
15     stratify=y # CANCELLARE QUESTA RIGA SE E' REGRESSIONE!
16 )
17
18 print(f"Shape Train: {X_train.shape}, Shape Test: {X_test.shape}")
```

Listing 8: Divisione Dati Corretta

## Clustering: KMeans & Ottimizzazione

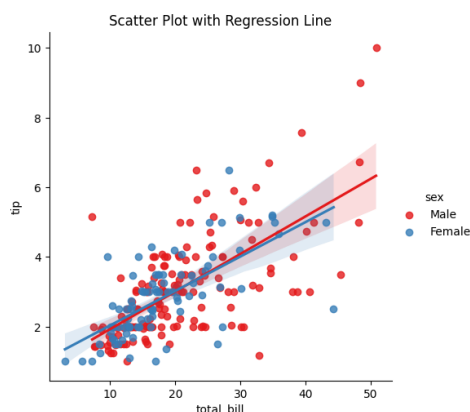
### 1. Concetti Chiave: Inerzia e Elbow Method

**Inerzia (Inertia):** È la somma delle distanze al quadrato tra ogni punto e il centro del suo cluster (centroide).

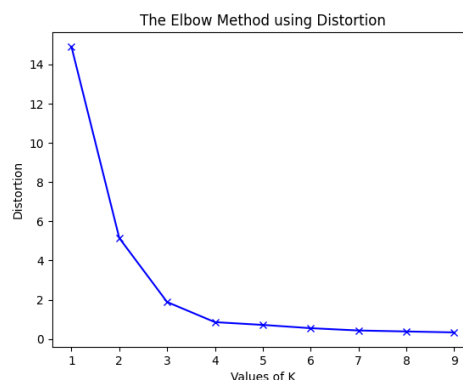
- Misura quanto i cluster sono "compatti".
- **Obiettivo:** Vogliamo un'inerzia bassa.
- **Il Problema:** Più aumentiamo  $K$  (numero di cluster), più l'inerzia scende (fino a 0 se  $K$  = numero di punti). Non possiamo solo minimizzarla, dobbiamo trovare un compromesso.

**Elbow Method (Metodo del Gomito):** Serve a scegliere il  $K$  ottimale. Grafichiamo l'inerzia al variare di  $K$ .

- Quando la curva scende ripida, aggiungere un cluster migliora molto il modello.
- Quando la curva si appiattisce (il "gomito"), aggiungere cluster non dà più grandi benefici ma complica solo il modello. Quel punto di flesso è il  $K$  ideale.



(a) Scatter Plot



(b) Elbow method, here,  $k = 3$

## 2. Codice Completo: Elbow Method e Scatter Plot

```
1
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import matplotlib.cm as cm
6 from sklearn.cluster import KMeans
7 from sklearn.metrics import silhouette_score, silhouette_samples
8 from sklearn.preprocessing import StandardScaler
9
10
11 # Se non hai X, commenta: X = df.select_dtypes(include=np.number).dropna()
12 scaler = StandardScaler()
13 X_scaled = scaler.fit_transform(X)
14
15 # --- PUNTO 4: K-MEANS OPTIMIZATION ---
16 print("--- 4. K-MEANS OPTIMIZATION ---")
17
18 # 1. Ricerca del miglior K (Constraints: k >= 3)
19 range_n_clusters = range(3, 11) # Da 3 a 10
20 best_k = -1
21 best_score = -1
22 best_model = None
23 best_labels = None
24
25 print(f"{'K':<5} | {'Silhouette Score'}")
26 print("-" * 25)
27
28 for k in range_n_clusters:
29     # Hyperparameters: n_init='auto' o 10 per stabilit , random_state fisso
30     kmeans = KMeans(n_clusters=k, init='k-means++', n_init=10, random_state=42)
31     cluster_labels = kmeans.fit_predict(X_scaled)
32
33     silhouette_avg = silhouette_score(X_scaled, cluster_labels)
34     print(f"{'k':<5} | {'silhouette_avg':.4f}")
35
36     if silhouette_avg > best_score:
37         best_score = silhouette_avg
38         best_k = k
39         best_model = kmeans
40         best_labels = cluster_labels
41
42 print("-" * 25)
43 print(f"Miglior schema trovato: K = {best_k} con Silhouette = {best_score:.4f}")
44
45 # 2. Mostrare gli Iperparametri (Richiesto esplicitamente)
46 print("\n--- Hyperparameters del modello migliore ---")
47 print(best_model.get_params())
48
49 # 3. Visualizzazione (Silhouette Plot + Distribuzione)
50 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))
51
52 # --- PLOT A: Silhouette Plot dei cluster ---
53 ax1.set_xlim([-0.1, 1])
54 ax1.set_ylim([0, len(X_scaled) + (best_k + 1) * 10])
55
56 sample_silhouette_values = silhouette_samples(X_scaled, best_labels)
57 y_lower = 10
58
59 for i in range(best_k):
60     ith_cluster_silhouette_values = sample_silhouette_values[best_labels == i]
61     ith_cluster_silhouette_values.sort()
62
63     size_cluster_i = ith_cluster_silhouette_values.shape[0]
64     y_upper = y_lower + size_cluster_i
65
66     color = cm.nipy_spectral(float(i) / best_k)
67     ax1.fill_betweenx(np.arange(y_lower, y_upper),
68                       0, ith_cluster_silhouette_values,
69                       facecolor=color, edgecolor=color, alpha=0.7)
70
71     ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))
72     y_lower = y_upper + 10
73
74 ax1.set_title(f"Silhouette Plot per i vari cluster (K={best_k})")
75 ax1.set_xlabel("Coefficiente di Silhouette")
76 ax1.set_ylabel("Etichetta Cluster")
77 ax1.axvline(x=best_score, color="red", linestyle="--", label="Media")
78 ax1.legend()
```

```

79 # --- PLOT B: Distribuzione delle label (Istogramma) ---
80 unique, counts = np.unique(best_labels, return_counts=True)
81 ax2.bar(unique, counts, color=cm.nipy_spectral(unique.astype(float) / best_k), alpha=0.7)
82 ax2.set_title(f"Distribuzione elementi per Cluster (K={best_k})")
83 ax2.set_xlabel("Cluster ID")
84 ax2.set_ylabel("Numero di Record")
85 ax2.set_xticks(unique)
86
87 # Aggiunge il conteggio sopra le barre
88 for i, v in enumerate(counts):
89     ax2.text(unique[i], v + 1, str(v), ha='center')
90
91 plt.tight_layout()
92 plt.show()
93

```

Listing 9: KMeans: Ricerca K ottimale e Visualizzazione

## Supervised Learning: Regressione, Classificazione e Tuning

In questa sezione affrontiamo l'addestramento dei modelli (con etichette, o *target*) e l'ottimizzazione degli iperparametri (GridSearch), tipici dei task finali dell'esame.

### 1. Il concetto di Grid Search e Cross-Validation

**Teoria:** Non possiamo sapere a priori quali parametri siano i migliori per un modello (es. quanto deve essere profondo un albero decisionale?).

- **GridSearch:** Prova "a forza bruta" tutte le combinazioni possibili di parametri che definiamo in una griglia.
- **Cross-Validation (CV):** Per evitare che il modello sia fortunato su un singolo split di dati, la CV divide il Train Set in  $K$  parti (fold). Allena su  $K - 1$  e valida sull'ultima. Ripete il processo  $K$  volte e fa la media dei punteggi.

### 2. Il "Mistero" del `neg_mean_squared_error`

Una domanda classica di teoria riguarda il parametro `scoring`.

**Perché usiamo il valore negativo dell'errore?**

- Le funzioni di ottimizzazione di `scikit-learn` (come `GridSearchCV`) sono programmate per **massimizzare** un punteggio (score). Cercano sempre il numero più alto possibile (es. Accuratezza 99% > 90%).
- L'errore (MSE o RMSE), invece, deve essere **minimizzato** (un errore di 0 è perfetto).
- **Soluzione:** Usiamo il negativo.

Errore 10  $\rightarrow$  Score  $-10$

Errore 2  $\rightarrow$  Score  $-2$

Poiché  $-2$  è matematicamente più grande (più a destra sulla retta dei numeri) di  $-10$ , l'algoritmo "massimizzando" il numero negativo sta in realtà minimizzando l'errore.

### 3. Codice Template: Regressione con GridSearch

Questo blocco è valido per i task che chiedono: *"optimize the depth... searching for minimum RMSE with cross-validation"*.

```

1 from sklearn.model_selection import GridSearchCV
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.metrics import mean_squared_error, r2_score
4 import numpy as np
5
6 # 1. Definizione del modello base (fissa il random_state!)
7 dt = DecisionTreeRegressor(random_state=42)

```

```

8
9 # 2. Definizione della griglia di parametri (Dizionario)
10 # Le chiavi devono corrispondere ESATTAMENTE ai parametri della funzione
11 param_grid = {
12     'max_depth': [3, 5, 10, 20, None],      # Profondità dell'albero
13     'min_samples_split': [2, 5, 10],        # Minimo nodi per dividere
14     'min_samples_leaf': [1, 2, 4]           # Minimo nodi in una foglia
15 }
16
17 # 3. Configurazione GridSearch
18 # cv=5: Divide i dati in 5 parti per testare la robustezza
19 # scoring='neg_mean_squared_error': Cerca di minimizzare l'errore
20 grid_search = GridSearchCV(
21     estimator=dt,
22     param_grid=param_grid,
23     cv=5,
24     scoring='neg_mean_squared_error',
25     n_jobs=-1, # Usa tutti i processori disponibili
26     verbose=1  # Mostra il progresso
27 )
28
29 # 4. Addestramento (Solo sul TRAIN set!)
30 print("Inizio ricerca parametri ottimali...")
31 grid_search.fit(X_train, y_train)
32
33 # 5. Recupero dei risultati migliori
34 print("\n--- Risultati Grid Search ---")
35 print("Migliori parametri trovati:", grid_search.best_params_)
36
37 # best_score_ il Negativo del MSE medio in Cross-Validation.
38 # Dobbiamo cambiarlo di segno e farne la radice per avere l'RMSE.
39 best_cv_rmse = np.sqrt(-grid_search.best_score_)
40 print(f"Miglior RMSE in Cross-Validation: {best_cv_rmse:.4f}")
41
42 # 6. Test sul Test Set (Valutazione Finale)
43 # grid_search.predict usa automaticamente il miglior modello trovato
44 y_pred = grid_search.predict(X_test)
45
46 final_mse = mean_squared_error(y_test, y_pred)
47 final_rmse = np.sqrt(final_mse)
48 print(f"RMSE finale sul Test Set: {final_rmse:.4f}")

```

Listing 10: GridSearch per DecisionTreeRegressor

## 4. Feature Selection basata sulla Correlazione

Spesso viene chiesto (es. Task 4 dell'esame) di rifare il modello scartando le colonne poco correlate.

```

1 # Calcolo correlazione con il target
2 target_col = 'price' # Sostituisci col nome vero del target
3 correlations = df.corr()[target_col].abs() # Valore assoluto!
4
5 # Soglia definita dall'esame (es. 0.15)
6 threshold = 0.15
7
8 # Seleziona le feature da tenere (quelle > 0.15)
9 selected_features = correlations[correlations > threshold].index.tolist()
10
11 # Rimuovi il target dalla lista delle feature X
12 if target_col in selected_features:
13     selected_features.remove(target_col)
14
15 print(f"Feature selezionate ({len(selected_features)}): {selected_features}")
16
17 # Creazione nuovi dataset ridotti
18 X_train_reduced = X_train[selected_features]
19 X_test_reduced = X_test[selected_features]
20
21 # Ora puoi ri-addestrare il modello su questi nuovi X_train_reduced

```

Listing 11: Rimuovere colonne con bassa correlazione

## 5. Classificazione e Metriche (Se richiesto)

Se l'esame chiede di predire una classe (es. "Sì/No" o "Tipo A/B/C") invece di un numero, cambia la metrica di scoring.



```

1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import accuracy_score, classification_report
3
4 # Nota: Classifier invece di Regressor
5 clf = DecisionTreeClassifier(random_state=42)
6
7 param_grid_clf = {
8     'criterion': ['gini', 'entropy'],
9     'max_depth': [3, 5, 10]
10 }
11
12 # Scoring cambia! Usiamo 'accuracy' o 'f1_macro'
13 grid_search_clf = GridSearchCV(
14     clf,
15     param_grid_clf,
16     cv=5,
17     scoring='accuracy' # 0 'f1_macro' per classi sbilanciate
18 )
19
20 grid_search_clf.fit(X_train, y_train)
21
22 print("Best Accuracy:", grid_search_clf.best_score_)
23 print("Best Params:", grid_search_clf.best_params_)

```

Listing 12: GridSearch per Classificazione

## DBSCAN (Density-Based Clustering)

A differenza del KMeans, il DBSCAN non richiede di specificare il numero di cluster ( $K$ ). Funziona in base alla **densità** dei punti.

### Quando usarlo?

- Quando i cluster hanno forme strane (non sferiche, es. lune, anelli).
- Quando ci sono molti *outliers* (rumore) che vogliamo ignorare.

### Parametri Critici:

- **eps** (Epsilon): Il raggio del cerchio intorno a ogni punto. Se è troppo piccolo, nessuno è vicino a nessuno (tutto rumore). Se è troppo grande, tutti sono collegati (un solo cluster gigante).
- **min\_samples**: Quanti vicini deve avere un punto per essere considerato un "Core Point" (parte centrale di un cluster).

```

1 from sklearn.cluster import DBSCAN
2
3 # Nota: DBSCAN molto sensibile allo Scaling! Assicurati di usare X_scaled.
4 # eps=0.5 e min_samples=5 sono valori di default, ma vanno quasi sempre cambiati.
5 dbscan = DBSCAN(eps=0.5, min_samples=5)
6 cluster_labels = dbscan.fit_predict(X_scaled)
7
8 # DBSCAN assegna etichetta -1 al "Rumore" (punti scartati)
9 n_clusters = len(set(cluster_labels)) - (1 if -1 in cluster_labels else 0)
10 n_noise = list(cluster_labels).count(-1)
11
12 print(f"Cluster trovati: {n_clusters}")
13 print(f"Punti di rumore (Noise): {n_noise}")
14
15 # Visualizzazione (Gestendo il colore del rumore)
16 plt.figure(figsize=(8, 6))
17 unique_labels = set(cluster_labels)
18 colors = [plt.cm.Spectral(each) for each in np.linspace(0, 1, len(unique_labels))]
19
20 for k, col in zip(unique_labels, colors):
21     if k == -1:
22         col = [0, 0, 0, 1] # Nero per il rumore
23
24     class_member_mask = (cluster_labels == k)
25     xy = X_scaled[class_member_mask]
26
27     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
28             markeredgecolor='k', markersize=6)

```

```

29
30 plt.title('DBSCAN Clustering')
31 plt.show()

```

Listing 13: DBSCAN Clustering

## Cheat Sheet: Interpretazione dei Risultati

Questa sezione serve per rispondere alle domande "Comment the results".

### 1. Metriche di Regressione (Numeri)

- **RMSE (Root Mean Squared Error):** È l'errore medio nella stessa unità di misura del target.  
*Esempio:* Se preveggo prezzi di case e RMSE = 50.000, sbaglio in media di 50k euro. **Più basso è, meglio è.**
- **R2 Score (R-Squared):** Spiega quanta varianza dei dati il modello ha catturato.  
Va da 0 a 1 (o negativo se il modello è disastroso).  
*Intervalli:* > 0.7 Buono, < 0.3 Scarso.

### 2. Metriche di Classificazione (Categorie)

- **Accuracy:** Percentuale di risposte esatte. Attenzione: se le classi sono sbilanciate (es. 99 sani, 1 malato), l'accuracy è ingannevole!
- **Confusion Matrix:**

$$\begin{bmatrix} VN & FP \\ FN & VP \end{bmatrix}$$

Vogliamo i numeri alti sulla **diagonale principale** (Veri Negativi e Veri Positivi). I numeri fuori diagonale sono gli errori.

### 3. Commenti Standard per l'Esame

*"The optimized Decision Tree performed better than the Linear Regression, showing a lower RMSE (12.5 vs 18.2). This suggests the relationship between features and target is non-linear."*

*"The DBSCAN found 3 clusters and identified 15 points as noise (outliers), which KMeans forced into the nearest cluster, reducing purity."*

## Modelli Avanzati: Random Forest

**Teoria:** Il Random Forest è un metodo "Ensemble". Invece di un solo albero (che rischia l'overfitting), ne costruisce tanti (es. 100) e fa votare ognuno. La risposta finale è la media (Regressione) o la maggioranza (Classificazione).

```

1 from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
2
3 # Esempio Classificatore
4 # n_estimators=100: Crea 100 alberi
5 rf = RandomForestClassifier(n_estimators=100, max_depth=10, random_state=42)
6
7 # Esempio Griglia per GridSearch con Random Forest
8 param_grid_rf = {
9     'n_estimators': [50, 100, 200], # Quanti alberi?
10    'max_depth': [5, 10, None],      # Quanto profondi?
11    'min_samples_split': [2, 5]
12 }
13
14 # Si usa nel GridSearch esattamente come il DecisionTree

```

Listing 14: Random Forest (Classifier o Regressor)

## Modello Alternativo: K-Nearest Neighbors (KNN)

**Teoria:** Il KNN non "impara" una formula. Per predire un nuovo punto, guarda i  $K$  punti più vicini nel training set.

- Se è Classificazione: prende la classe di maggioranza dei vicini.
- Se è Regressione: fa la media dei valori dei vicini.

**Nota:** Richiede assolutamente lo **Scaling** dei dati!

```
1 from sklearn.neighbors import KNeighborsClassifier
2
3 # n_neighbors il parametro 'K'
4 knn = KNeighborsClassifier(n_neighbors=5)
5
6 param_grid_knn = {
7     'n_neighbors': [3, 5, 7, 9], # Dispari per evitare pareggi nel voto
8     'weights': ['uniform', 'distance'] # 'distance' d  pi  peso ai vicini stretti
9 }
```

Listing 15: KNN Classifier

## Interpretazione del Modello: Feature Importance

Spesso nel punto "Comment the results" è utile dire quale variabile ha influito di più. Questo codice funziona solo con Decision Trees e Random Forest.

```
1 # Assumiamo che 'best_model' sia il tuo albero/foresta gi  addestrato
2 import pandas as pd
3
4 # Ottieni l'importanza
5 importances = best_model.feature_importances_
6 feature_names = X_train.columns # I nomi delle colonne
7
8 # Crea un DataFrame per ordinare i dati
9 feat_imp_df = pd.DataFrame({'Feature': feature_names, 'Importance': importances})
10 feat_imp_df = feat_imp_df.sort_values(by='Importance', ascending=False)
11
12 # Grafico a barre
13 plt.figure(figsize=(10, 6))
14 sns.barplot(x='Importance', y='Feature', data=feat_imp_df, palette='viridis')
15 plt.title("Quali feature contano di pi  ?")
16 plt.show()
```

Listing 16: Grafico Feature Importance

**Commento d'Esame:** *"From the Feature Importance plot, we see that 'Income' is the most dominant predictor (0.6 importance), while 'Age' has almost no impact on the decision. This suggests focusing business logic on Income."*