: # Importa la libreria NumPy co import numpy as np # Importa la funzione train_te	st_split dal modulo model_selection della libreria scikit-learn
np.random.seed(0) # Genera 100 campioni casuali altezze = np.random.normal(0, # Genera i pesi corrispondenti pesi = 0.5 * altezze + np.rand # Suddivide le altezze e i pes # Il 30% dei dati è assegnato # random_state=42 garantisce l X_train, X_test, y_train, y_te # Stampa le dimensioni del set print("Dimensioni del Training # Stampa le dimensioni del set	per la riproducibilità dei risultati da una distribuzione normale con media 0 e deviazione standard 5, rappresentanti le altezze 5, 100) alle altezze, dove il peso è calcolato come metà dell'altezza più un termine di errore casuale
Dimensioni del Training Set (a Dimensioni del Test Set (altez 2 - Analisi della Rela 2 - Analisi della Rela 3 - Analisi della Rela 3 - Analisi della Rela 4 - Analisi della Rela 5 - Analisi della Rela 5 - Analisi della Rela 6 - Analisi della 7 - Analisi della Rela 6 - Analisi della 7 -	ltezze e pesi): (70,) (70,) ze e pesi): (30,) (30,) azione tra Visite al Sito e Importo delle Vendite ie NumPy per la generazione di dati casuali t # Importa Matplotlib per la visualizzazione grafica mport train_test_split # Importa train_test_split per la suddivisione del dataset la riproducibilità dei risultati
<pre>visite_al_sito = np.random.ran importo_vendite = 50 + 0.2 * v # Suddivide il dataset in set X_train, X_test, y_train, y_te # Crea un grafico a dispersion plt.figure(figsize=(10, 6)) # plt.scatter(X_train, y_train, plt.scatter(X_test, y_test, la plt.xlabel('Numero di Visite a plt.ylabel('Importo delle Vendo</pre>	iglia
	Set (visite al sito e importo delle vendite):", X_train.shape, y_train.shape) (visite al sito e importo delle vendite):", X_test.shape, y_test.shape) Relazione tra Visite al Sito e Importo delle Vendite
100 - 150 - 150 - 100 -	
Dimensioni del Test Set (visit	Numero di Visite al Sito isite al sito e importo delle vendite): (700,) (700,) e al sito e importo delle vendite): (300,) (300,) porzioni delle Classi nei Set di Addestramento e Test
<pre>import numpy as np # Importa np.random.seed(1) # Imposta if # Genera 100 campioni casuali X = np.random.rand(100, 2) # Genera 100 etichette casuali Y = np.random.choice(['A', 'B']) # Calcola la proporzione della proporzione_classe_A = sum(Y =</pre>	<pre>classe 'A' nel vettore delle etichette = 'A') / len(Y)</pre>
<pre># Suddivide il dataset in set X_train, X_test, y_train, y_te # Calcola la proporzione della proporzione_classe_A_train = s # Calcola la proporzione della proporzione_classe_B_train = 1 # Calcola la proporzione della</pre>	di addestramento (70%) e set di test (30%) st = train_test_split(X, Y, test_size=0.3, random_state=42) classe 'A' nel set di addestramento um(y_train == 'A') / len(y_train) classe 'B' nel set di addestramento - proporzione_classe_A_train
# Calcola la proporzione della proporzione_classe_B_test = 1 # Stampa le proporzioni calcol print("Proporzione Classe A ne print("Proporzione Classe B ne print("Proporzione Classe B ne print("Proporzione Classe B ne print("Proporzione Classe A ne print("Proporzione Classe B ne print("Proporzione Classe B ne print("Proporzione Classe B ne	classe 'B' nel set di test - proporzione_classe_A_test ate 1 data Set completo:", proporzione_classe_A) 1 data Set completo:", proporzione_classe_B) 1 Training Set:", proporzione_classe_A_train) 1 Training Set:", proporzione_classe_B_train) 1 Test Set:", proporzione_classe_A_test) 1 Test Set:", proporzione_classe_B_test) Set completo: 0.54 Set completo: 0.54 Set completo: 0.459999999999996
Proporzione Classe B nel Train Proporzione Classe A nel Test Proporzione Classe B nel Test 4-Visualizzazione (import matplotlib.pyplot as pl # Definizione delle etichette labels = ['Classe A', 'Classe # Definizione dei colori per l colors = ['gold', 'lightcoral'	set: 0.5666666666666666666666666666666666666
	lassi nel Set')
54.0%	
5-Analisi Statistica : import random # Importa il mo	di un Campione Casuale dulo random per generare numeri casuali il modulo NumPy per calcoli numerici a lista vuota per il dataset
<pre>for i in range(1000): dataset.append(random.rand) # Estrazione di un campione ca campione_casuale = random.samp # Calcolo della media e della media_campione = np.mean(campi deviazione_standard_campione = # Calcolo della media e della media_dataset = np.mean(datase)</pre>	int(1, 100)) # Aggiunge un numero casuale compreso tra 1 e 100 al dataset suale semplice di 300 elementi dal dataset le(dataset, 300) # Estrae 300 elementi casuali dal dataset deviazione standard del campione one_casuale) # Calcola la media del campione np.std(campione_casuale) # Calcola la deviazione standard del campione deviazione standard del dataset completo t) # Calcola la media del dataset completo np.std(dataset) # Calcola la deviazione standard del dataset completo np.std(dataset) # Calcola la deviazione standard del dataset completo
# Stampa dei risultati print(f"Media del campione cas print(f"Deviazione standard de print(f"Media del dataset comp print(f"Deviazione standard de Media del campione casuale: 50 Deviazione standard del campio Media del dataset completo: 50 Deviazione standard del datase	uale: {media_campione:.2f}") 1 campione casuale: {deviazione_standard_campione:.2f}") 1 campione casuale: {deviazione_standard_campione:.2f}") 1 dataset completo: {deviazione_standard_dataset:.2f}") 1.57 1.57 1.60 1.7
<pre>import numpy as np # Per la import matplotlib.pyplot as pl</pre>	manipolazione dei dati con DataFrame generazione di numeri casuali t # Per la visualizzazione dei dati mport train_test_split # Per la suddivisione del dataset roducibilità dei risultati
# Generare la colonna con dist	<pre>', 'B'], size=num_elementi, p=[percentuale_A, 1 - percentuale_A]) : colonna})</pre>
1 B 2 B 3 A 4 A 995 A 996 B 997 A 998 B	
999 A 1000 rows × 1 columns 7-Creazione di Tre : # Creare tre subset di dimensi	Subset con Dimensioni Simili dal DataFrame oni simili le del DataFrame df con una frazione del 1/3
<pre># Rimuovi dal DataFrame df gli df = df.drop(subset1.index) # Estrai un sottoinsieme casua subset2 = df.sample(frac=1/2) # Rimuovi dal DataFrame df gli df = df.drop(subset2.index)</pre>	
8 2500000 50000 125 20 8-Calcolo delle Per # Calcolare le percentuali di	rcentuali di 'A' e 'B' nel Subset1 "A" e "B" per il subset1 'ColonnaAB'].value_counts (normalize=True)
<pre># Calcolare le percentuali di percentuali_subset1 = subset1[percentuali_subset2 = subset2[percentuali_subset3 = subset3[# Creare i grafici a torta fig, axs = plt.subplots(3, 1, # Plot per Subset 1 axs[0].pie(percentuali_subset1]</pre>	Centuali di 'A' e 'B' per Ogni Subset "A" e "B" per ogni subset 'ColonnaAB'].value_counts(normalize=True) # Calcola le percentuali per il subset1 'ColonnaAB'].value_counts(normalize=True) # Calcola le percentuali per il subset2 'ColonnaAB'].value_counts(normalize=True) # Calcola le percentuali per il subset3 figsize=(6, 12)) # Crea una griglia di subplot 3x1 con dimensioni della figura (6, 12) , labels=percentuali_subset1.index, autopct='%1.1f%%', startangle=90) # Crea un grafico a torta per il Subset 1
<pre>axs[0].set_title('Subset 1') # Plot per Subset 2 axs[1].pie(percentuali_subset2 axs[1].set_title('Subset 2') # Plot per Subset 3 axs[2].pie(percentuali_subset3 axs[2].set_title('Subset 3') # Mostrare il grafico plt.show() # Mostra tutti i s</pre>	# Imposta il titolo del subplot 1 , labels=percentuali_subset2.index, autopct='%1.1f%%', startangle=90) # Crea un grafico a torta per il Subset 2 # Imposta il titolo del subplot 2 , labels=percentuali_subset3.index, autopct='%1.1f%%', startangle=90) # Crea un grafico a torta per il Subset 3 # Imposta il titolo del subplot 3
70.6% A	В
Subset 2 30.5%	В
69.5% A Subset 3	
73.6% A	
<pre>: # Dividere ciascun subset in t train_subset1, test_subset1 = train_subset2, test_subset2 =</pre>	train_test_split(subset1, test_size=0.2, random_state=42) train_test_split(subset2, test_size=0.2, random_state=42) train_test_split(subset3, test_size=0.2, random_state=42) e
<pre>def draw_pie(ax, data, title): ax.pie(data, labels=data.i ax.set_title(title) # Prima riga di torte (Subset draw_pie(axs[0, 0], train_subset draw_pie(axs[0, 1], test_subset # Seconda riga di torte (Subset draw_pie(axs[1, 0], train_subset draw_pie(axs[1, 1], test_subset # Terza riga di torte (Subset)</pre>	ndex, autopct='%1.1f%%', startangle=90) 1) et1['ColonnaAB'].value_counts(normalize=True), 'Train Subset 1') t1['ColonnaAB'].value_counts(normalize=True), 'Test Subset 1') t 2) et2['ColonnaAB'].value_counts(normalize=True), 'Train Subset 2') t2['ColonnaAB'].value_counts(normalize=True), 'Test Subset 2')
71.1% A	B 31.3% 68.7% A
Train Subset	B 28.5%
71.5% A	A 61.2%
Train Subset	B 4.4%
75.6% A 11-Analisi Statistic	a di un DataFrame
<pre>data = {'Valori': [1, 2, 3, 4, df = pd.DataFrame(data) # Calcola la media dei valori mean_value = df['Valori'].mean_value</pre>	con una colonna chiamata "Valori" 5, 10, 15, 20, 25, 300, 1000, 100000000, -50000000, -50]} nella colonna "Valori" () rd dei valori nella colonna "Valori"
: # Identifica gli outliers cons	a dei Valori in un DataFrame iderando ±3 sigma dalla media mean_value + 3 * std_dev) (df['Valori'] < mean_value - 3 * std_dev)]
: # Crea un grafico a dispersion plt.scatter(df.index, df['Valcout # Evidenzia gli outlier nel graplt.scatter(outliers.index, outliers.index, outliers.index)	ri'], label='Valori') afico con un colore rosso tliers['Valori'], color='red', label='Outliers') le per rappresentare la media
<pre># Aggiunge linee orizzontali p plt.axhline(y=mean value + 3 *</pre>	
1e8 Grafi 1.00 - Valori Outliers 0.75 - Media ±3 Deviazioni \$ 0.50 - 0.25 - 0.25 - 0.00	co con Outliers Evidenziati
-0.25 - -0.50 - -0.75 -	4 6 8 10 12 Indice
<pre>import pandas as pd # Importa import matplotlib.pyplot as pl # Crea un DataFrame di esempio data = {'Feature1': [1, 2000,</pre>	degli Outliers nelle Features del DataFrame la libreria pandas per la manipolazione dei dati t # Importa la libreria matplotlib per la visualizzazione dei dati con 4 features 3, 4, 50000, 10, 15, 20, 2500000, 300000000, 100000000], 8, 10, 20, 30, 40, 50000, 60, 200], 1, 20000, 25, 50, 75, 100, 125, 150, 500000], 1, 20000, 3, 4000000000, 5, 10, 15, 20, 20005, 30, 10000]}
<pre># Definisci il numero minimo o min_features_threshold = 1 k = 2 # Fattore di moltiplica # Lista per salvare gli indici outlier_indices = [] # Itera su ogni feature for feature in df.columns: mean_value = df[feature].m std_dev = df[feature].std # Identifica gli outliers</pre>	ean()) per ciascuna feature e li salva come una nuova colonna nel DataFrame
<pre>df['Outlier_' + feature] = # Visualizza il DataFrame con df Feature1 Feature2 Feature3 0 1 2 5</pre>	(df[feature] > mean_value + k * std_dev) (df[feature] < mean_value - k * std_dev) le colonne aggiuntive per gli outlier Feature4
5 10 20 50 6 15 30 75 7 20 40 100 8 2500000 50000 125 9 300000000 60 150 10 100000000 200 500000	5 False False False False False 10 False False False False 15 False False False False 20 False False False False 2005 False True False False 30 True False False False 10000 False True False False True False False True False False True False
<pre># Calcola il numero totale di df['Num_Outliers'] = df.filter df Feature1 Feature2 Feature3 0 1 2 5 1 2000 4 10 - 2 3 6 15</pre>	outlier per ogni riga nel DataFrame df (like='Outlier_').sum(axis=1) Feature4 Outlier_Feature1 Outlier_Feature2 Outlier_Feature3 Outlier_Feature4 Num_Outliers 1 False False False False False O 20000000 False False False False False O 3 False False False False False O
3 4 8 20000 40 4 50000 10 25 5 10 20 50 6 15 30 75 7 20 40 100 8 2500000 50000 125 9 300000000 60 150 10 100000000 200 500000	False False False True 1 False False False True 1 False False False False False True 1 False False False False False 0 False False False False False 0 False False False False False 0 False False False False False 1 True False False False 1
<pre>: # Calcola il numero di feature df['Num_Outliers'] = df.filter # Filtra i dati per mantenere outliers = df[df['Num_Outliers # Aggiungi una colonna che ind df['Is_Outlier'] = df.index.is # Rimuovi colonne ausiliarie df.drop(df.filter(like='Outlier')</pre>	solo le righe con almeno il numero minimo di features superanti la soglia '] >= min_features_threshold] ica se il record è un outlier o meno in (outliers.index) r_').columns, axis=1, inplace=True)
0 1 2 5 1 2000 4 10 - 2 3 6 15 3 4 8 20000 40 4 50000 10 25	Feature Is_Outlier 1 False 2000000 False 3 False 0000000 True 5 False
5 10 20 50 6 15 30 75 7 20 40 100 8 2500000 50000 125 9 300000000 60 150 10 100000000 200 500000	10 False 15 False 20 False 2005 True 30 True 10000 True de dei Grafici in una Matrice
# Organizza i grafici in una mum_features = len(df.columns) num_features #num_features vie # se il DataFrame ha 5 colonne 4 18-Visualizzazione # Definizione del numero di ri num_rows = num_features	atrice, con una colonna e 4 righe - 1 # Escludi la colonna 'Is_Outlier' ne assegnata al numero di colonne nel DataFrame meno uno, poiché stiamo escludendo la colonna 'Is_Outlier' , di cui una è 'Is_Outlier', allora num_features sarà 4, poiché stiamo contando solo le colonne che rappresentano le feature effettive. delle Features con Outliers Evidenziati ghe e colonne per la visualizzazione dei dati
<pre>num_rows = num_features num_cols = 1 # Una colonna # Creazione di una nuova figur plt.figure(figsize=(6, 4 * num # Iterazione attraverso le col for i, feature in enumerate(df</pre>	a matplotlib con dimensioni specificate _rows)) onne del DataFrame per visualizzare i dati .columns[:-1]): # Escludi la colonna 'Is_Outlier' afico nella posizione specificata cols, i + 1) ter plot dei dati eature], c=df['Is_Outlier'], cmap='coolwarm', alpha=0.8) del sotto-grafico sso - {feature}') x
<pre># Etichettatura dell'asse plt.ylabel('Feature') # Ottimizzazione del layout de plt.tight_layout() # Visualizzazione della figura plt.show()</pre> 1e8 00 2.5 -	$i\ sotto-grafici$
2.0 - 9.1.5 - 1.0 - 0.5 - 0.0 - 2	
0 2 OI 50000 - 40000 -	4 6 8 10 Indice Itliers in Rosso - Feature2
20000 - 10000 - 0 2	4 6 8 10 Indice
500000 - 400000 - 300000 -	
T	itliers in Rosso - Feature3 •
100000 -	
100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000	attliers in Rosso - Feature3
100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 10000	elle Righe Corrispondenti agli Outliers
100000 - 100000 - 1e9 Ot 1e9 Ot 3.5 - 3.0 - 2.5 - 1.0 - 0.5 - 0.0 - 0 z 19-Eliminazione del DataFrame df df_filtered = df[df['Is_Outlies # Visualizzazione del DataFrame df_filtered]	ettiers in Rosso - Feature3 titiers in Rosso - Feature4 titiers in Rosso - Feature4 ettiers in Rosso
100000 - 100000 - 100000 - 1e9 Otto 4.0 - 3.5 - 3.0 - 2.5 - 1.0 - 0.5 - 0.0 - 2 19-Eliminazione del DataFrame del del filtered = del [del ['Is_Outlies # Visualizzazione del DataFrame del filtered 1 2 5 1 2000	etitiers in Rosso - Feature3
100000 le9 4.0 3.5 3.0 2.5 1.0 0.5 0.0 2 19-Eliminazione de # Elimina le righe corrisponde # Filtraggio del DataFrame df df filtered = df [df ['Is_Outlie # Visualizzazione del DataFram df_filtered 10 1 2 5 1 2000 4 10 -2000 2 3 6 15 4 50000 10 25 5 10 20 50 6 15 30 75 7 20 40 100 20-Calcolo della D def calcola deviazione_standar # Calcola la media della i media = sum (lista) / # Calcola la somma dei que somma_quadrati_diff = sum # Calcola la deviazione standar # Calcola la deviazione standar # Calcola la somma dei que somma_quadrati_diff = sum # Calcola la deviazione standar # Esempio di utilizzo nueviazione_standard = (som return deviazione_standard # Esempio di utilizzo nueviazione_standard = (som return deviazione_standard # Esempio di utilizzo nueviazione_standard = (som return deviazione_standard # Esempio di utilizzo nueviazione_standard = (som return deviazione_standard # Esempio di utilizzo nueviazione_standard = (som return deviazione_standard # Esempio di utilizzo nueviazione_standard = (som return deviazione_standard # Esempio di utilizzo nueviazione_standard # Esempio di utilizzo	stillers in Rosso - Feature 3 **Black Rights Corrispondenti agli Outliers **Black Rights Corrispondenti agli Outliers **Black Rights Corrispondenti agli Outliers **To Add and the Rosso - Feature 4 **To Add and the Rosso - Feature 5 **To Add