

CHIARA BERETTA
1033576

k-means

Ai e machine learning per il marketing

Obiettivi del progetto:

Comprendere e sperimentare il funzionamento degli algoritmi di machine learning attraverso la creazione di un modello di segmentazione clienti con k-means attraverso l'utilizzo dell'intelligenza artificiale

Strumenti utilizzati:

- ChatGPT e Cloude
- JupyterLab
- Streamlit
- GitHub

Perchè il k-medie ?

Il k-medie è un algoritmo non supervisionato che suddivide i dati in gruppi omogenei in base alla distanza da un centroide, rendendolo particolarmente adatto alla segmentazione dei clienti nel marketing.

Per rappresentare questa logica, ho scelto di lavorare su un esempio vicino alla realtà: un e-commerce dedicato al pet care e alla cura degli animali.

Ogni cluster riflette un diverso profilo d'acquisto e l'algoritmo consente di assegnare dinamicamente nuovi utenti al gruppo più coerente con le loro abitudini, simulando il comportamento di un sito che utilizza i dati raccolti per proporre strategie su misura

Come funziona il k-medie ?

L'algoritmo delle k-medie raggruppa un insieme di dati in k gruppi distinti (cluster), basandosi sulla vicinanza tra i punti.

I passaggi:

- Si scelgono il numero di cluster (k) da creare
- Vengono inizializzati casualmente k centroidi (punti centrali di ciascun cluster)
- Ogni punto del dataset viene assegnato al centroide più vicino
- Si calcolano nuovi centroidi per ogni gruppo usando la media dei punti assegnati.
- I punti vengono riassegnati ai nuovi centroidi.
- Il processo si ripete finché i centroidi non cambiano più

Alla fine, ogni punto è associato al cluster il cui centroide è più vicino

Fase 1: creazione del database con Chatgpt e Cloude

PROMT:

"Crea un esempio in python di segmentazione clienti usando l'algoritmo k-means. I dati devono rappresentare 4 gruppi di clienti, con caratteristiche sintetiche ma realistiche: spesa media mensile e numero di acquisti. Visualizza i cluster ottenuti e assegna un nuovo cliente al gruppo più vicino. Genera un codice eseguibile su JupyterLab, con una visualizzazione grafica chiara che evidenzi anche i centroidi e il nuovo cliente."

Codice creato:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
```

Impostazione della librerie che permettono di gestire dati numerici, visualizzarli graficamente ed eseguire la segmentazione tramite l'algoritmo K-Means

```
np.random.seed(42)

cluster1 = np.column_stack((
    np.random.normal(loc=160, scale=6, size=50),
    np.random.normal(loc=4, scale=1.2, size=50)
))

cluster2 = np.column_stack((
    np.random.normal(loc=100, scale=6, size=50),
    np.random.normal(loc=10, scale=1.2, size=50)
))

cluster3 = np.column_stack((
    np.random.normal(loc=130, scale=6, size=50),
    np.random.normal(loc=5, scale=1.2, size=50)
))

cluster4 = np.column_stack((
    np.random.normal(loc=115, scale=6, size=50),
    np.random.normal(loc=13, scale=1.2, size=50)
))

X = np.vstack((cluster1, cluster2, cluster3, cluster4))
X[:, 1] = np.clip(X[:, 1], a_min=1, a_max=None)

nuovo_cliente = np.array([[120, 6]])
```

Generazione dei dati: simulazione di clienti con caratteristiche diverse

```
kmeans = KMeans(n_clusters=4, random_state=0)
kmeans.fit(X)
```

Addestramento del modello: applicazione dell'algoritmo K-Means

```
labels = kmeans.labels_
centroids = kmeans.cluster_centers_
cluster_assegnato = kmeans.predict(nuovo_cliente)

print(f"Il nuovo cliente {nuovo_cliente[0]} è stato assegnato al cluster:", cluster_assegnato[0])
```

Predizione: assegnazione del cluster al nuovo cliente

```
plt.figure(figsize=(10, 7))
plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='Accent', edgecolor='k', s=60, alpha=0.75)
plt.scatter(centroids[:, 0], centroids[:, 1], c='red', marker='X', s=250, label='Centroidi')
plt.scatter(nuovo_cliente[0, 0], nuovo_cliente[0, 1],
            facecolors='white', edgecolors='gray',
            marker='*', s=300, linewidths=3, label='Nuovo Cliente')
plt.xlabel("Spesa media mensile (€)")
plt.ylabel("Numero di acquisti mensili")
plt.title("Segmentazione clienti con K-Means")
plt.legend()
plt.grid(True)
plt.tight_layout()
plt.show()
```

Visualizzazione: rappresentazione grafica dei cluster, dei centroidi e del nuovo punto

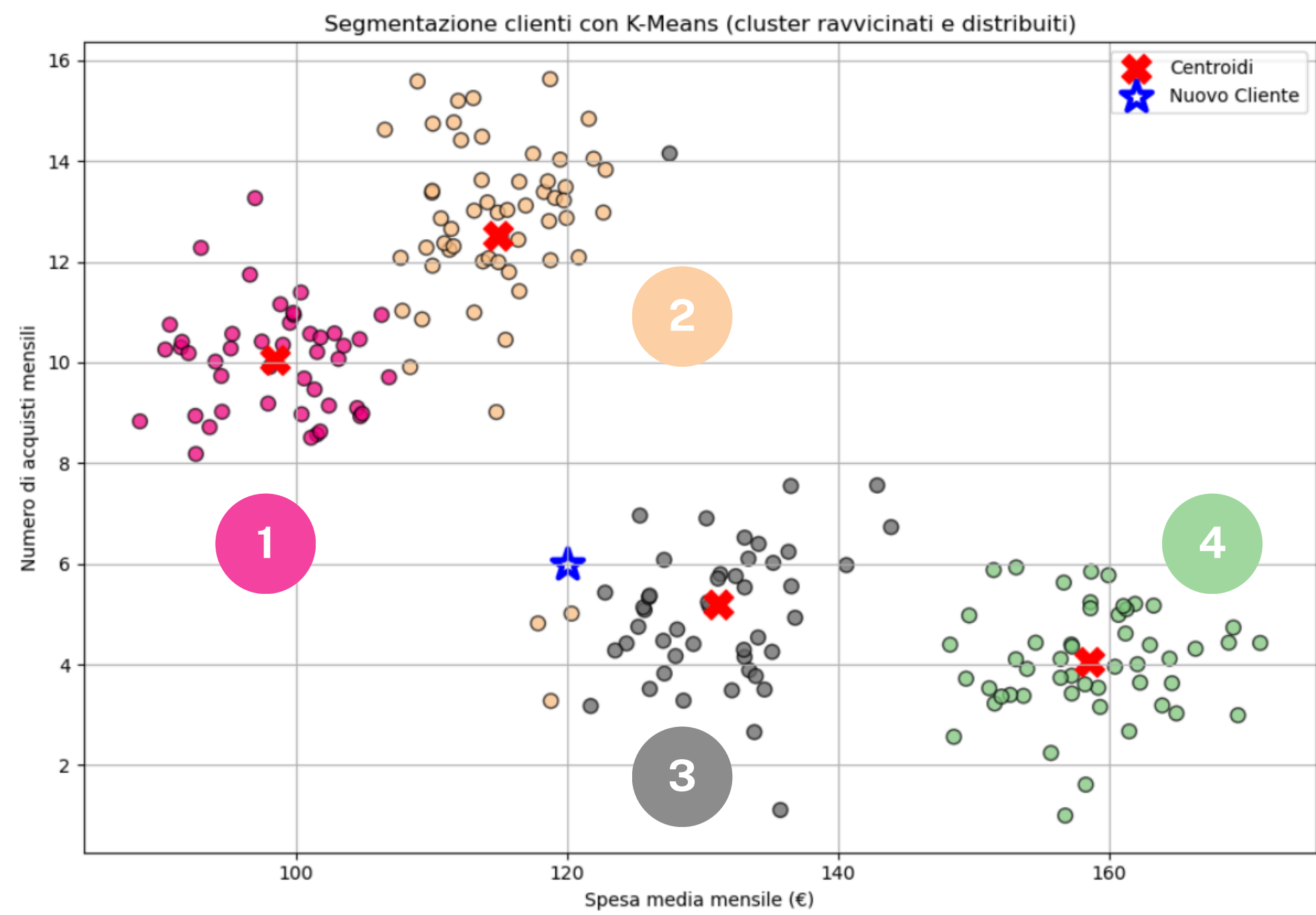
Campione:

Formato da 200 osservazione

Variabili:

- X = spesa media mensile (€)
- Y = numero di acquisti mensili

Visualizzazione:



Visualizzazioe con JupyterLab

Descrizione cluster individuati:

1 Cacciatori di offerte

Spesa media mensile: €100 circa

Numero di acquisti mensili: 10 + σ

Descrizione: clienti che effettuano acquisti frequenti ma di poco valore (ex. puntando a offerte e promozioni su prodotti essenziali come cibo secco, lettiera e snack). Sono fidelizzati perchè fanno molto acquisti ma molto attenti al risparmio

3 Animali affezionati

Spesa media mensile: €130 circa

Numero di acquisti mensili: 5 + σ

Descrizione: clienti con un comportamento regolare, che acquistano cibo, prodotti e per la cura. Rappresentano la base stabile e prevedibile dell'e-commerce.

2 Pet lover organizzati

Spesa media mensile: €115 circa

Numero di acquisti mensili: 13 + σ

Descrizione: clienti molto attivi che effettuano acquisti molto spesso. Alterneranno quindi prodotti base ad articoli extra (ex. giochi, snack funzionali, antiparassitari), pianificando gli ordini con attenzione.

4 Amici animali premium

Spesa media mensile: €160 circa

Numero di acquisti mensili: 4 + σ

Descrizione: clienti che acquistano meno frequentemente ma scelgono prodotti di alta qualità (ex. crocchette biologiche, accessori premium, integratori specifici ...)

Strategie per ogni cluster:

1 Cacciatori di offerte

- Newsletter settimanale con sconti mirati
- Bundle promozionali (ex. 2+1 su snack)
- Programma fedeltà a punti → premia la quantità di acquisti
- Sezione “offerte della settimana” in homepage

2 Pet lover organizzati

- Abbonamenti personalizzati
- Reminder per riacquisto (ex. email o notifica: "Sta per finire la scorta di...")
- Promo early-access: “prenota ora le novità prima degli altri”

3 Animali affezionati

- Email mensile personalizzata con consigli
- Sconti alle ricorrenze (ex. anniversario primo ordine, compleanno pet ...)
- Up-selling mirato
- Test nuovi prodotti gratuiti

4 Amici animali premium

- Upselling mirato
- Packaging premium e personalizzato
- Storytelling dei prodotti via e-mail
- Programma VIP esclusivo

Un nuovo cliente

Il nuovo cliente creato casualmente da ChatGPT ha assunto delle variabili casuali.

In base a questi dati, l'algoritmo k-means ha assegnato il cliente al cluster 2, denominato "Pet lover organizzati". Questo comportamento dopo un po' di mesi di osservazione suggerisce che il cliente ha già iniziato ad acquistare con regolarità sul nostro sito dimostrando coinvolto e soddisfazione della propria esperienza. Quindi anche su di esso si posso attivare le strategie relative al cluster 2.

Fase 2: quale strategia marketing è più adatta per il singolo?

Dal modello statico all'app interattiva

Dopo la simulazione e visualizzazione dei cluster con JupyterLab, l'obiettivo è stato quello di simulare una reale applicazione dell'algoritmo k-means con altri dati esterni al database inseriti casualmente su un'app web

Perché? questa applicazione potrebbe simulare il reale funzionamento di un sistema automatico di clusterizzazione e raccomandazione quando si raccolgono abbastanza dati su un nuovo cliente

Dalla staticità all'interattività

Basandomi sul database scritto in Python in precedenza si sono presi i centroidi di ogni cluster insieme alla deviazione standard di ognuno

Dopo ho utilizzato Streamlit per trasformare il modello in un'app:



- Teoricamente il sito dopo un determinato periodo di utilizzo di un nuovo utente elabora due valori (spesa e acquisti) ma all'interno dell'app web lo faremo noi
- Analizzando questi due valori con il modello k-means lo si assegna al cluster più vicino
- L'app restituisce:
 - il profilo assegnato
 - le strategie consigliate
 - il grafico dinamico con il punto evidenziato per verificare la coerenza


Pubblicazione online

Strumenti utilizzati per la messa online dell'app: GitHub + Streamlit Cloud

Per condividere pubblicamente l'app ho seguito questi passaggi:

- Caricato il codice su GitHub (app.py + requirements.txt)
- Collegato il repository a Streamlit Cloud
- Lanciato il sito web disponibile a tutti:

 <https://app-petcare-5zytjelay.pr4p2glklgeql.streamlit.app/> 



Pubblicazione online: app.py

```
1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 # Palette colori e descrizioni
7 palette = {
8     "Cacciatori di offerte": "deeppink",
9     "Pet lover organizzati": "sandybrown",
10    "Animali affezionati": "dimgray",
11    "Amici animali premium": "mediumseagreen"
12 }
13
14 strategie = {
15     "Cacciatori di offerte": [
16         📧 "Newsletter settimanale",
17         🍷 "Bundle 2+1 snack",
18         🏆 "Programma punti",
19         🛒 "Offerte in homepage"
20     ],
21     "Pet lover organizzati": [
22         🐾 "Abbonamenti personalizzati",
23         🕒 "Reminder riacquisto",
24         ⚡ "Promo early-access"
25     ],
26     "Animali affezionati": [
27         📧 "Email mensile personalizzata",
28         🎁 "Sconti ricorrenze",
29         💎 "Up-selling mirato",
30         🌱 "Test nuovi prodotti"
31     ],
32     "Amici animali premium": [
33         💎 "Upselling mirato",
34         🍷 "Packaging premium",
35         📧 "Storytelling via email",
36         🏆 "Programma VIP"
37     ]
38 }
39
40 # Simulazione dati coerente con il tuo grafico
41 np.random.seed(42)
42 cluster_specs = {
43     "Cacciatori di offerte": {"mean": (100, 10), "std": (5, 1), "count": 50},
44     "Pet lover organizzati": {"mean": (115, 13), "std": (5, 1), "count": 50},
45     "Animali affezionati": {"mean": (130, 5), "std": (5, 1), "count": 50},
46     "Amici animali premium": {"mean": (160, 4), "std": (5, 1), "count": 50},
47 }
48
49 data = []
50 for cluster, specs in cluster_specs.items():
51     x = np.random.normal(specs["mean"][0], specs["std"][0], specs["count"])
52     y = np.random.normal(specs["mean"][1], specs["std"][1], specs["count"])
53     for xi, yi in zip(x, y):
54         data.append({
55             "Spesa": xi,
56             "Acquisti": yi,
57             "Cluster": cluster
```

```
53     for xi, yi in zip(x, y):
54         data.append({
55             "Spesa": xi,
56             "Acquisti": yi,
57             "Cluster": cluster
58         })
59
60 df = pd.DataFrame(data)
61
62 # Centroidi stimati
63 centroids = {
64     "Cacciatori di offerte": (100, 10),
65     "Pet lover organizzati": (115, 13),
66     "Animali affezionati": (130, 5),
67     "Amici animali premium": (160, 4)
68 }
69
70 # Funzione per assegnare cluster
71 def assegna_cluster(x, y):
72     distanze = {k: np.sqrt((x - cx) ** 2 + (y - cy) ** 2) for k, (cx, cy) in centroids.items()}
73     return min(distanze, key=distanze.get)
74
75 # Streamlit UI
76 st.title("🐾 Scopri il tuo profilo cliente PetCare!")
77 st.markdown("""
78 Inserisci la **spesa media mensile** e il **numero di acquisti** per scoprire a quale cluster appartieni e quali strategie sono più adatte a te.
79 """)
80
81 spesa = st.number_input("\U0001F4B0 Spesa media mensile (€)", min_value=0.0, value=100.0, step=1.0)
82 acquisti = st.number_input("\U0001F6CD Numero di acquisti mensili", min_value=0, value=5, step=1)
83
84 profilo = assegna_cluster(spesa, acquisti)
85 st.markdown(f"### 🐾 Il tuo cluster è: **{profilo}**")
86 st.markdown(f"**\U0001F4CC Strategie consigliate per questo profilo:**")
87 for s in strategie[profilo]:
88     st.write(f"- {s}")
89
90 # Grafico
91 fig, ax = plt.subplots(figsize=(10, 6))
92 for cluster in df["Cluster"].unique():
93     subset = df[df["Cluster"] == cluster]
94     ax.scatter(subset["Spesa"], subset["Acquisti"], label=cluster, color=palette[cluster], edgecolors='black')
95
96 for cluster, (x, y) in centroids.items():
97     ax.scatter(x, y, color=palette[cluster], marker='X', s=200, edgecolors='black', label=None)
98
99 # Nuovo punto utente
100 ax.scatter(spesa, acquisti, color="gray", marker='*', s=200, edgecolors='black', label="Nuovo Cliente")
101
102 ax.set_xlabel("Spesa media mensile (€)")
103 ax.set_ylabel("Numero di acquisti mensili")
104 ax.set_title("Segmentazione clienti con K-Means")
105 ax.legend()
106 st.pyplot(fig)
```

Codice Python per la
visualizzazione dell'applicazione

Pubblicazione online: requirements.txt

Installazione delle librerie:

- 1 `scikit-learn` → Per eseguire il clustering con l'algoritmo k-means
- 2 `streamlit` → Per costruire l'interfaccia interattiva dell'app web
- 3 `numpy` → Per creare e manipolare i dati simulati
- 4 `matplotlib` → Per visualizzare i cluster e i centroidi su grafico

Senza questo file, l'app non potrebbe funzionare correttamente online

Cosa ho imparato

- Ho compreso passo per passo la logica delle k-medie, approfondendo come si formano i cluster, cosa rappresentano i centroidi e in che modo un nuovo punto viene automaticamente assegnato al gruppo più vicino
- Ho generato un database realistico e visualizzato i gruppi con grafici chiari, grazie a un utilizzo consapevole dell'intelligenza artificiale generativa e delle librerie Python per la simulazione e l'analisi dei dati
- Ho sviluppato una web app interattiva con Streamlit e GitHub, trasformando il modello in un'interfaccia semplice, accessibile da browser e in grado di simulare ciò che farebbe un e-commerce reale dopo aver raccolto i dati dei clienti.

Grazie per l'attenzione !