

Programmazione

Gli strumenti per i Big Data rappresentano l'infrastruttura fondamentale che consente di gestire, analizzare e visualizzare enormi quantità di dati eterogenei provenienti da fonti diverse.

Questi strumenti si dividono in più categorie: piattaforme di storage e gestione dei dati (come Hadoop Distributed File System e Apache HBase), strumenti di elaborazione e calcolo distribuito (Apache Spark, Flink, MapReduce) e soluzioni di integrazione e streaming (come Apache Kafka e NiFi).

Il loro scopo principale è garantire la scalabilità orizzontale cioè la capacità di distribuire i dati e i processi su più nodi e la tolleranza ai guasti, elementi chiave per l'affidabilità dei sistemi Big Data.



Parallelamente, esistono strumenti orientati all'analisi, alla visualizzazione e al machine learning, che permettono di estrarre informazioni utili dai dati raccolti.

Framework come TensorFlow, PyTorch o Scikit-learn vengono spesso integrati con i sistemi Big Data per eseguire analisi predittive o classificazioni su larga scala.

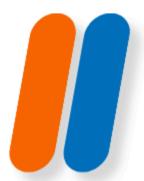
Inoltre, strumenti di data visualization come Tableau, Power BI o Apache Superset consentono di trasformare i risultati delle analisi in dashboard interattive e comprensibili, rendendo i dati realmente accessibili ai decisori aziendali.



In sintesi, l'ecosistema dei Big Data si basa su un insieme di tecnologie complementari che lavorano in sinergia, dall'acquisizione fino alla rappresentazione dei risultati.

#### Storage e gestione dei dati

- Hadoop HDFS (Hadoop Distributed File System) → Sistema di file distribuito che consente di memorizzare grandi volumi di dati su cluster di macchine, garantendo scalabilità e tolleranza ai guasti.
- Apache HBase → Database NoSQL basato su colonne, progettato per l'accesso rapido a dati non strutturati e per integrarsi con Hadoop.
- MongoDB → Database NoSQL orientato ai documenti, ideale per gestire dati semi-strutturati come JSON in modo flessibile e scalabile.



#### Storage & gestione dati

Teoria (HDFS / MongoDB).

Nel livello storage i requisiti chiave sono scalabilità orizzontale, replica e throughput. HDFS distribuisce file in blocchi su più nodi con replica per tolleranza ai guasti; è ottimo per dati "append-only" e workload batch. I database NoSQL (es. MongoDB) offrono schemi flessibili e accesso a bassa latenza su dati semi-strutturati. Spesso si usa HDFS come data lake "grezzo" e un NoSQL per le query operative.

```
1.# Carica un CSV locale su HDFS (cartella /data/raw)
2.hdfs dfs -mkdir -p /data/raw
3.hdfs dfs -put customers.csv /data/raw/
4.
5.# Lista e lettura
6.hdfs dfs -ls /data/raw
7.hdfs dfs -cat /data/raw/customers.csv | head -n 5
```

```
1.from pymongo import MongoClient
2.
3.# Connessione a MongoDB (default locale)
4.client = MongoClient("mongodb://localhost:27017")
5.db = client["retail"]
6.coll = db["customers"]
7.
8.# Inserimento documento (JSON-like)
9.doc = {"id": 1001, "name": "Alice", "city": "Torino", "spend": 320.5}
10.coll.insert_one(doc)
11.
12.# Query con filtro e proiezione
13.for c in coll.find({"city": "Torino"}, {"_id": 0, "name": 1, "spend": 1}):
14. print(c)
```



#### Elaborazione e calcolo distribuito

- Apache Spark → Framework per il calcolo distribuito in memoria, molto più veloce di MapReduce, usato per analisi batch e realtime.
- Apache Flink → Sistema per l'elaborazione di flussi di dati in tempo reale (stream processing), utile per applicazioni che richiedono aggiornamenti continui.
- MapReduce → Modello di programmazione che suddivide grandi volumi di dati in blocchi elaborati in parallelo, poi aggregati per ottenere risultati finali.



#### Elaborazione e calcolo distribuito

Teoria (Spark).

I motori come Apache Spark eseguono trasformazioni su RDD/DataFrame distribuendo i task sul cluster. Il calcolo in memoria riduce la latenza rispetto a MapReduce. Il modello è "lazy": le trasformazioni costruiscono un piano logico, le azioni lo materializzano. Spark è adatto a batch, ML e, con Structured Streaming, anche a near real-time.

```
1.from pyspark.sql import SparkSession
2.from pyspark.sql.functions import explode, split, col
3.
4.spark = SparkSession.builder.appName("WordCount").getOrCreate()
5.
6.# Carico un testo su HDFS o locale (qui: locale per semplicità)
7.df = spark.read.text("alice.txt") # una riga per record
8.
9.# Tokenizzo e conteggio parole
10.words = df.select(explode(split(col("value"), r"\W+")).alias("word")).where(col("word")!= "")
11.counts = words.groupBy("word").count().orderBy(col("count").desc())
12.
13.counts.show(10, truncate=False)
14.spark.stop()
```



#### Integrazione e streaming dei dati

- Apache Kafka → Piattaforma di streaming distribuita che gestisce grandi flussi di dati in tempo reale, spesso usata per connettere sistemi diversi.
- Apache NiFi → Strumento per l'automazione dei flussi di dati (dataflow) tra sistemi, con un'interfaccia grafica per definire pipeline di elaborazione.
- Airflow → Framework per la gestione e l'automazione dei workflow di dati, utile per pianificare e monitorare processi ETL complessi.



#### Integrazione e streaming dei dati

Teoria (Kafka).

Apache Kafka è un log distribuito publish/subscribe a throughput elevato. I producer scrivono su topic partizionati; i consumer leggono in gruppi, scalando orizzontalmente.

È usato per ingest in tempo reale e come "backbone" eventi tra microservizi. La durabilità su disco e la retention configurabile lo rendono adatto a re-processing.

```
1.# Producer
                                                                    1.# Consumer
                                                                    2.from kafka import KafkaConsumer
2.from kafka import KafkaProducer
                                                                    3. import json
3. import json
                                                                    5.consumer = KafkaConsumer(
5. producer = KafkaProducer(
                                                                       "events_purchases",
    bootstrap_servers=["localhost:9092"],
                                                                       bootstrap_servers=["localhost:9092"],
    value_serializer=lambda v: json.dumps(v).encode("utf-8")
                                                                       auto_offset_reset="earliest",
8.)
                                                                       value_deserializer=lambda v: json.loads(v.decode("utf-8")),
                                                                       group_id="analytics-g1"
10. for i in range(5):
                                                                   11.)
    event = {"event_id": i, "type": "purchase", "amount": 9.99 + i}
    producer.send("events_purchases", value=event)
                                                                   13. for msg in consumer:
13.
                                                                   14. print("EVENT:", msg.value)
14. producer.flush()
                                                                   15. break # demo: leggi un messaggio e termina
```



## Analisi, machine learning e intelligenza artificiale

- TensorFlow → Libreria open source di Google per il deep learning e l'analisi su larga scala.
- PyTorch → Framework di machine learning sviluppato da Meta, noto per la flessibilità nello sviluppo di reti neurali.
- Scikit-learn → Libreria Python per l'apprendimento automatico tradizionale (classificazione, regressione, clustering) integrabile nei flussi Big Data.



# Analisi, machine learning e intelligenza artificiale

Teoria (scikit-learn / Spark ML).

Per ML "tradizionale" su dataset medio-grandi si usa spesso scikitlearn (feature engineering, crossvalidation, pipeline).

Per scala elevata si passa a Spark MLlib, che distribuisce sia trasformazioni sia training.

Best practice: Pipeline end-to-end, train/test split, metriche riproducibili e versionamento dei modelli.

```
1. import pandas as pd
 2.from sklearn.model_selection import train_test_split
 3. from sklearn.preprocessing import StandardScaler
 4. from sklearn.compose import ColumnTransformer
 5. from sklearn.pipeline import Pipeline
 6.from sklearn.linear_model import LogisticRegression
 7.from sklearn.metrics import classification_report
 9.# Dataset di esempio
10.df = pd.DataFrame({
     "age": [25, 45, 36, 52, 23, 40],
     "income": [30_000, 60_000, 45_000, 80_000, 28_000, 55_000],
    "city": ["TO", "MI", "BO", "MI", "TO", "BO"],
14. "label": [0, 1, 0, 1, 0, 1]
15.})
16.
17.X = df[["age", "income", "city"]]
18.y = df["label"]
20.num_cols = ["age", "income"]
21.cat_cols = ["city"]
23.pre = ColumnTransformer([
24. ("num", StandardScaler(), num_cols),
25. # One-hot encoding automatico con get_dummies pre-fit: per semplicità uso pandas
26.], remainder="passthrough")
28.X_enc = pd.get_dummies(X, columns=cat_cols) # semplice OHE
30. Xtr, Xte, ytr, yte = train_test_split(X_enc, y, test_size=0.33, random_state=42)
32.pipe = Pipeline([
    ("scaler", StandardScaler(with_mean=False)), # with_mean=False per sparse
    ("clf", LogisticRegression(max_iter=500))
35.])
36.
37. pipe.fit(Xtr, ytr)
38.pred = pipe.predict(Xte)
39.print(classification_report(yte, pred))
```

#### Visualizzazione e business intelligence

- Tableau Software di visualizzazione dati che consente di creare dashboard interattive e comprensibili anche a utenti non tecnici.
- Power BI → Strumento Microsoft per analizzare e rappresentare dati aziendali integrandosi con Excel, SQL Server e servizi cloud.
- Apache Superset → Piattaforma open source per la visualizzazione e l'esplorazione di dataset Big Data, integrabile con database SQL e NoSQL.



## Visualizzazione e business intelligence

Teoria (Superset / Power BI / Tableau).

I tool BI si collegano a data warehouse/lakehouse (es. Presto/Trino, DuckDB, BigQuery, Redshift, Spark SQL), eseguono SQL e forniscono dashboard interattive con filtri, controlli e condivisione. La logica resta lato DB; il tool gestisce caching, permessi, e governance di dashboard. La chiave è mantenere modelli semantici chiari (metriche, dimensioni, calcoli) e query performanti (indici, partizionamento).

```
1.-- Aggregazione vendite per città e fascia di età
2.SELECT
3. city,
4. CASE
5. WHEN age < 30 THEN 'U30'
6. WHEN age BETWEEN 30 AND 49 THEN '30-49'
7. ELSE '50+'
8. END AS age_band.
9. COUNT(*) AS customers,
13. WHERE dt BETWEEN DATE '2025-09-01' AND DATE '2025-09-30'
```

