

Fabric in un caso reale di analisi sinistri con open AI

Marco Pozzan



Sponsor



Chi sono

Consulente e formatore in ambito business intelligence

Nel 2023 Fondatore e CTO start-up Regolo Farm per Business Intelligence e Power Platform

Nel 2020 Fondatore e CTO start-up Cargo BI per insurance data analytics per le compagnie assicurative

Dal 2017 mi occupo di architetture big Data e in generale di tutta la proposizione data platform di Microsoft.

Docente all'Università di Pordenone per i corsi IFTS di analisi Big Data

Community Lead di 1nn0va (www.innovazionefvg.net)

MCP,MCSA,MCSE dal 2017 MCT e dal 2014 MVP per SQL Server e relatore in diverse conferenze sul tema.

- Marco.pozzan@regolofarm.com
- [@marcopozzan.it](https://www.marcopozzan.it)
- www.marcopozzan.it



Microsoft
CERTIFIED
Trainer



Microsoft
CERTIFIED
Professional
Microsoft
CERTIFIED
Solutions Associate
SQL Server 2012/2014

Agenda

- Che cosa è Fabric?
- Introduzione alla Demo
 - Che cosa è LLM?
 - Che cosa sono i RAG?
 - Cosa è LangChain?
 - Collegamento a Fabric da Notebook

WHAT IS MICROSOFT FABRIC?

A UNIFIED SAAS
PLATFORM FOR ALL YOUR
ANALYTICS NEEDS

REAL-TIME ANALYTICS

INGEST & QUERY REAL-TIME DATA
WITHIN SECONDS

STORE REALTIME DATA IN OneLake

BI

REPORTS &
DASHBOARDS
WITH
POWER BI

DATA SCIENCE

RUN NOTEBOOKS WITHIN SECONDS
WITHOUT MANAGING CLUSTERS
TRAIN MODELS AND EXPERIMENT
BETWEEN DIFFERENT MODELS VERSIONS

DATA WAREHOUSE

POWERED BY DELTA PARQUET
ACCESSIBLE THROUGH SQL AND SPARK
DECOUPLED STORAGE & PROCESSING
ORGANIZED AROUND WORKSPACES

OneLake

UNIFIED
LAKEHOUSE
WITH AN OPEN
FORMAT

DATA INTEGRATION

USE NO-CODE OR NOTEBOOKS FOR
ETL OR ELT JOBS
SCHEDULE & RUN JOBS WITHOUT
MANAGING (SPARK) CLUSTERS



UNIFIED SECURITY
ACROSS ALL TYPES OF
DATA STORES & ENGINES

SHORTCUTS EASY ACCESS TO DATA STORED IN OTHER CLOUDS

STORE DATA



LAKEHOUSE



WAREHOUSE



KQL DATABASE



NOTEBOOK



DATAFLOW GEN2



DATA PIPELINE



KQL QUERYSET



SPARK JOB
DEFINITION



EXPERIMENT



MODEL

TRAIN AND USE ML MODELS



SCORECARD



REPORT



DASHBOARD



PAGINATED
REPORT



REALTIME
DASHBOARD

FABRIC
ITEMS

ALL-IN-ONE

PRESENT DATA

MANAGE REAL-TIME DATA



EVENTSTREAM



STREAMING
DATAFLOW

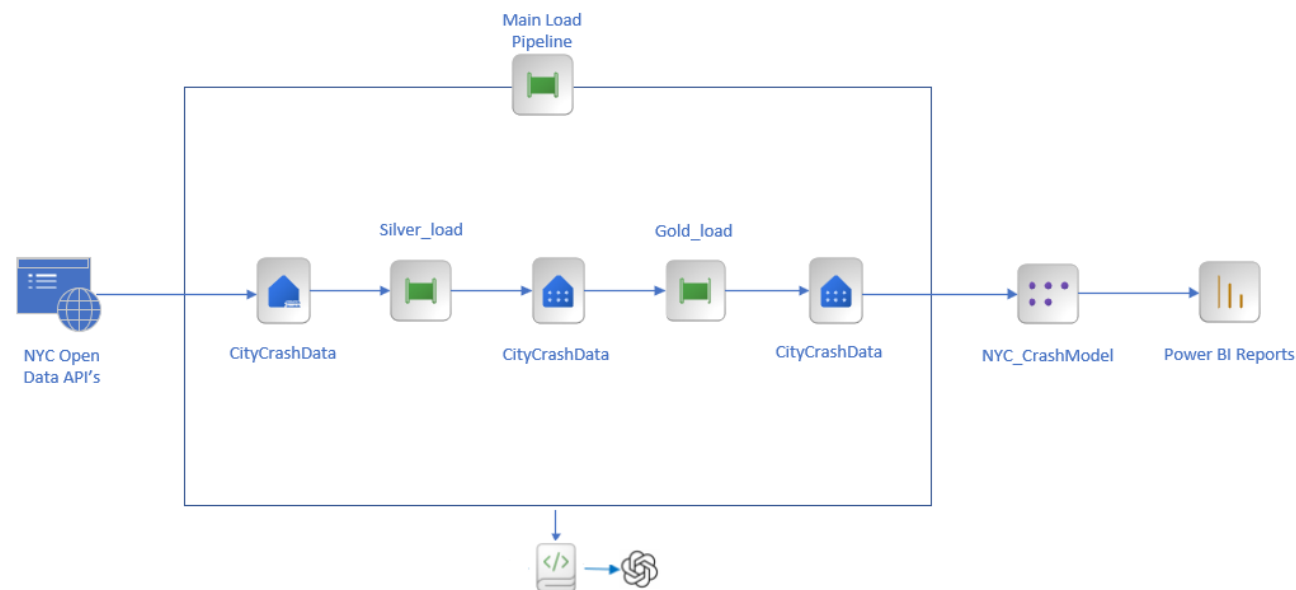


STREAMING
DATASET

Demo

A New York City, negli Stati Uniti, è richiesto un rapporto di polizia standard quando si verifica una collisione di veicoli in cui qualcuno rimane ferito, muore o se ci sono almeno \$ 1.000 di danni. Questi dati vengono raccolti e pubblicati online come fonte di dati aperta su NYC Open Data e vengono aggiornati quotidianamente.

Lo scenario è quello di fornire report Power BI aggiornati quotidianamente per fornire analisi descrittive e approfondimenti su posizione, fattori contribuenti e tipi di veicoli per collisioni di veicoli utilizzando Fabric. Questa soluzione fornisce anche le basi per altre persone come Data Analyst e Data Scientist per analisi con Open AI.



Che cosa è LLM?

- I Large Language Model (LLM), come ChatGPT di OpenAI, hanno rivoluzionato il modo in cui interagiamo e utilizziamo l'intelligenza artificiale.
- Questi modelli sono abili nel generare testo simile a quello umano e possono essere sfruttati per un'ampia gamma di applicazioni. Tuttavia, per massimizzare la loro utilità in contesti specifici, come ambienti aziendali o domini specializzati, la personalizzazione è essenziale.

Che cosa è la Retrieval Augmented Generation (RAG)?

- Di base, gli LLM sono formati su ampi set di dati che comprendono un'ampia gamma di conoscenze generali. Sebbene ciò li renda versatili, significa anche che potrebbero non avere la precisione necessaria per attività specializzate.
- Un metodo efficace per migliorare la pertinenza e l'accuratezza degli output LLM è tramite Retrieval Augmented Generation (RAG). RAG integra le risposte del modello con informazioni recuperate dinamicamente, assicurando che l'output sia contestualmente rilevante e aggiornato. Ciò è particolarmente prezioso quando si ha a che fare con dati aziendali specifici o query dei clienti, in cui le risposte generiche risultano carenti.

Che cosa è LangChain

- LangChain è un framework per lo sviluppo di applicazioni basate su modelli linguistici di grandi dimensioni (LLM).
- LangChain semplifica ogni fase del ciclo di vita dell'applicazione LLM:
- **Sviluppo:** crea le tue applicazioni utilizzando i blocchi di costruzione open source, i componenti e le integrazioni di terze parti di LangChain. Utilizza LangGraph per creare agenti con stato con streaming di prima classe e supporto umano nel ciclo.
- **Produzione:** utilizza LangSmith per ispezionare, monitorare e valutare le tue catene, in modo da poterle ottimizzare e distribuire in modo continuo con sicurezza.
- **Distribuzione:** trasforma le tue applicazioni LangGraph in API e assistenti pronti per la produzione con LangGraph Cloud.

Come si configura LangChain

- Codice che inizializza l'istanza `chat_model` di LangChain che fornisce un'interfaccia per richiamare un provider LLM tramite chat API.
- Il motivo per cui selezioniamo il **modello chat** è che il modello **gpt-35-turbo** è ottimizzato per la chat

```
1  import os
2  from langchain.chat_models import AzureChatOpenAI
3
4  llm = AzureChatOpenAI(model=os.getenv("OPENAI_CHAT_MODEL"),
5                        deployment_name=os.getenv("OPENAI_CHAT_MODEL"),
6                        temperature=0)
```

- La temperatura è un parametro che controlla la "creatività" o la casualità del testo generato. Una temperatura più bassa rende l'output più "focalizzato" o deterministico. Poiché qui abbiamo a che fare con un database, è importante che LLM generi una risposta relativa ai fatti.

Configurazione del prompt per LLM

- **Prompt** è l'input che inviamo all'LLM per generare un output. Prompt può anche essere progettato per contenere istruzioni, contesto e impostare il tono e formattare i dati di output.

```
1  from langchain.prompts.chat import ChatPromptTemplate
2
3  final_prompt = ChatPromptTemplate.from_messages(
4      [
5          ("system",
6           """
7           You are a helpful AI assistant expert in querying SQL Database to find answers to user's question about Collision.
8           """,
9          ),
10         ("user", "{question}\n ai: "),
11     ]
12 )
```

- Utilizzare Prompt Template è un buon modo per strutturare queste proprietà, incluso l'input dell'utente finale da fornire all'LLM.
- Noi utilizziamo il modulo **ChatPromptTemplate** di LangChain, che si basa su **ChatML** (Chat Markup Language).

Configurazione dell'agent

- Il componente **Agent** di LangChain è un wrapper attorno a LLM, che decide i migliori passaggi o azioni da intraprendere per risolvere un problema.
- L'Agent ha in genere accesso a un set di funzioni chiamate **Tools** (o Toolkit) e può decidere quale Tool utilizzare in base all'input dell'utente.
- Ogni agent può eseguire varie attività NLP, come parsing, calcoli, traduzione ecc.

Configurazione dell'agent

- Un **Agent Executor** è un'interfaccia eseguibile dell'Agent e del suo set di Tools.
- L'agent executor è responsabile della chiamata dell'Agent, dell'ottenimento dell'azione e dell'input dell'azione, della chiamata dello strumento a cui l'azione fa riferimento con l'input corrispondente, dell'ottenimento dell'output dello strumento e quindi del passaggio di tutte queste informazioni all'Agent per ottenere l'azione successiva che dovrebbe intraprendere. Di solito è un processo iterativo finché l'Agent non raggiunge la Final Answer o l'output.

Come ci si connette al lakehouse di Fabric

- Per collegarsi al lakehouse di fabric usiamo la libreria sqlalchemy passando come credenziali ClientID e Secret del Service Principal

```
1  from sqlalchemy import create_engine
2  import sqlalchemy as sa
3
4  driver = '{ODBC Driver 18 for SQL Server}'
5  odbc_str = sa.URL.create(
6      "mssql+pyodbc",
7      username="feca4d00-63fe-47b6-95d1-4fe2ca30dd48",
8      password="tsd8Q~GoCYMTvUGYeJamAUtgBxku.6ovPZAeVbQy",
9      host="ehauyipajrjepfj3id6yqalcnm-rxx6zuqhqlterk4ko7sdyhvdsm.datawarehouse.fabric.microsoft.com",
10     database="CityCrashData",
11     query={"driver": "ODBC Driver 18 for SQL Server", "Authentication": "ActiveDirectoryServicePrincipal"}
12 )
13 print(odbc_str)
14 db_engine = create_engine(odbc_str)
15 connection = db_engine.raw_connection()
```

Configurazione dell'agent

- L'agent è progettato per interagire con il database SQL. E' dotato di toolkit per connettersi al database SQL e leggere sia i metadati che dati

```
1  from langchain.agents import AgentType, create_sql_agent
2  from langchain.sql_database import SQLDatabase
3  from langchain.agents.agent_toolkits.sql.toolkit import SQLDatabaseToolkit
4
5  db = SQLDatabase(db_engine, view_support=True)
6  # test the connection
7  print(db.dialect)
8  print(db.get_usable_table_names())
9
10 sql_toolkit = SQLDatabaseToolkit(db=db, llm=llm)
11 sql_toolkit.get_tools()
12
13 sqldb_agent = create_sql_agent(
14     llm=llm,
15     toolkit=sql_toolkit,
16     agent_type=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
17     verbose=True
18 )
```

- Abbiamo utilizzato ZERO_SHOT_REACT_DESCRIPTION come valore del parametro agent_type, che indica all'agente di non utilizzare memoria.

Come funziona l'agent?

- L'agent è progettato per mandare domande e ricevere risposte.

```
1  sqldb_agent.run(final_prompt.format(  
2      question="what is total of persone morte"  
3  ))  
4  # question="what is the number of persone morte in the year 2023?"
```

Grazie!