

Architectures for big data - Assignment 1

Alessandro Di Gioacchino, Yousef Hammar, Mattia Paravisi

November 9, 2022

1 Introduzione

Questo assignment ha come obiettivo quello di costruire una architettura usando le classi astratte di Python. L'architettura dovrebbe aiutare un possibile team di developer a raggiungere il seguente business requirement: "I need to show Intercompany impacts on my Company Balance Sheet, without any impact on OneStream performance during the Month End Closing activities".

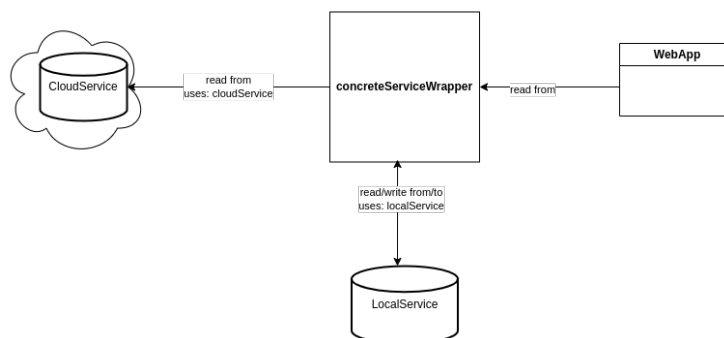
2 Premesse

Durante la risoluzione dell'assignment nel nostro gruppo sono state fatte delle premesse per limitare i possibili scenari che avremmo dovuto affrontare:

1. Abbiamo un servizio online che possiede un database dal quale dobbiamo leggere.
2. Il servizio online ha almeno una tabella di log da cui vogliamo leggere.
3. Dal servizio online ci limitiamo a leggere con lo scopo di creare un mirror locale.
4. Ogni tabella di log ha un timestamp.
5. Abbiamo un database locale.
6. Il database locale ha una tabella che viene usata per fare il mirroring della tabella online.
7. Dal database locale si può leggere e sul database locale si può scrivere.
8. Non ci interessa restituire alla web application i dati più aggiornati

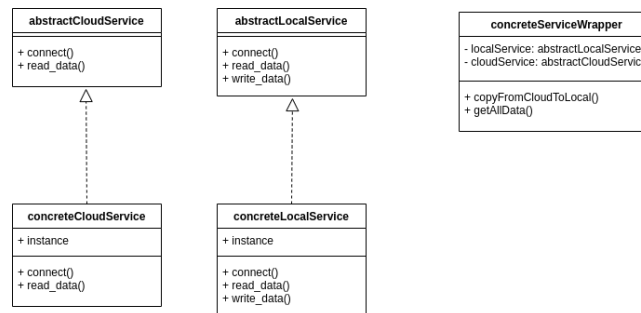
3 Architettura

Abbiamo riassunto la situazione come presentato nel seguente schema:



Consideriamo un servizio online - che deve avere un database associato - da cui vogliamo leggere, un database locale su cui possiamo scrivere e dal quale possiamo leggere e un differente servizio online (ad esempio una web application) che vuole accedere ai dati di cui stiamo facendo il mirroring. Vogliamo fare il

mirroring dei dati per non intaccare sulle performance delle procedure che il primo servizio esegue a fine mese. La nostra architettura consiste in un software che esegue tutte le operazioni richieste utilizzando come supporto delle classi astratte che verranno adattate ad ogni servizio utilizzato quando rese concrete. In particolare abbiamo immaginato la seguente gerarchia:



avremo quindi una classe astratta per il servizio online e una classe astratta per il servizio in locale. Le due classi astratte sono molto simili tra loro ma abbiamo deciso di mantenere le due implementazioni separate per permettere di avere per la classe che si occupa del servizio cloud solo il metodo read, mentre alla classe che si occupa del servizio locale sia read che write. Così facendo, inoltre, possiamo forzare la classe wrapper a prendere come argomenti i due tipi nel modo corretto; avendo un'unica gerarchia e non le due classi separate avremmo potuto rischiare di invertire localService e cloudService a causa di un errore. Per quanto riguarda la classe ServiceWrapper abbiamo deciso di non fare altro che considerare i due metodi utili al fine di raggiungere i business requirements:

- Il primo metodo copia i dati dal servizio cloud al servizio locale. (Fa il mirroring)
- Il secondo metodo restituisce tutti i dati da una certa tabella.

3.1 Esempio di codice

Uno snippet di codice molto ad alto livello che utilizza questo codice è il seguente:

```

1 localService = concreteLocalService(...)
2 cloudService = concreteCloudService(...)
3 serviceWrapper = concreteServiceWrapper(localService, cloudService)
4 serviceWrapper.copyFromCloudToLocal(...)

```

nella chiamata serviceWrapper.copyFromCloudToLocal(...):

```

1 def copyFromCloudToLocal(...):
2     res = cloudService.read_data(...)
3     localService.write_data(res)

```

Così facendo le classi concrete per il servizio locale e cloud possono essere qualsiasi; non ci interessa su che dbms si basano i due servizi, ci basta implementare un metodo read e un metodo write coerente alla tecnologia utilizzata. Dei parametri possibili per il metodo read (e di conseguenza per write) potrebbero essere:

```
1 def copyFromCloudToLocal(...):  
2     res = cloudService.read_data(table, timestamp)  
3     localService.write_data(table, res)
```

in questo modo permetteremmo la lettura delle righe di una specifica tabella che hanno un timestamp consono e la scrittura nella corrispondente tabella locale.

Consideriamo l'implementazione della seguente classe:

```

1  select = {
2      "exists": "",
3      "columns": [],
4      "from": "",
5      ...
6      "where": "",
7      "groupBy": [],
8      ...
9  }
10
11  class QueryBuilder():
12
13      def __init__(self):
14          self.query = ""
15          self.selectdict = select.copy()
16
17      def columns(self, columns):
18          self.selectdict["columns"] = columns
19
20      def from_table(self, table):
21          self.selectdict["from"] = table
22
23      def where(self, selectdict: dict, where):
24          self.selectdict["where"] = where
25
26      def build_select(self):
27          self.query = ""
28          if len(self.selectdict["columns"]) > 0:
29              self.query = "SELECT " + ', '.join(self.selectdict
30              ["columns"])
31          ...
32          return self.query

```

Come è facile intuire la classe descritta permette di creare delle query qualsiasi, potremmo quindi utilizzarla per costruire le query che verranno utilizzate nei metodi "read_data()" e "write_data()" delle due classi descritte sopra, ad esempio:

```

1  def read_data(table, timestamp):
2      builder = QueryBuilder()
3      builder.from_table("prova")
4      builder.where("timestamp_column > " + timestamp)
5      self.engine.execute(builder.build_select())

```

in cui engine è ottenuto utilizzando il metodo connect nel seguente modo:

```

1  def connect(self, url):
2      self.engine = create_engine(url)
3      return

```

usando sqlalchemy ad esempio, url può essere una stringa del tipo:

- "mysql://root:password@127.0.0.1:3306/public"
- "postgresql://root:password@127.0.0.1:5432/public"

Come si nota non ci interessa il dbms utilizzato, ci basta riuscire a connetterci.

4 Connessione ai pillars

1. Being the framework for satisfying requirements: in questo caso il requirement era uno solo come riportato nell'introduzione; una volta definita

l'architettura è stato possibile verificare se il requirement fosse stato raggiunto: la nostra architettura permette di fare il mirroring del servizio durante i giorni in cui non vengono effettuate le operazioni di fine mese, quindi effettivamente non inficiamo sulle performance del servizio.

2. Being the managerial basis for cost estimation and process management: avendo un'architettura decisa permette di eseguire un'analisi dei costi. In questo caso si dovrebbero valutare i costi OPEX come il costo del servizio cloud, i costi degli sviluppatori e dei servizi locali. Inoltre sappiamo quante persone impiegare per sviluppare il numero di classi descritto nell'architettura.
3. Enabling component reuse: la nostra architettura permette di fare il mirroring da due generici servizi, quindi è altamanete riutilizzabile.
4. Avoiding handover and people lock-in: l'architettura è semplice da documentare in quanto le classi che vengono utilizzate sono poche come i metodi che ogni classe deve implementare.