



FLIGHT DATA ANALYSIS

SECONDO PROGETTO BIG DATA



Alessandro Wood & Luca Gregori

Contesto dei Dati



Dati relativi ai voli domestici
negli USA dal 2009 al 2019.

Raccolti dal Dipartimento
dei Trasporti americano.

Un Dataset per ogni anno.



TECNOLOGIE UTILIZZATE



Spark, Spark Streaming, Spark MLlib

Kafka

HDFS Hadoop

Cassandra

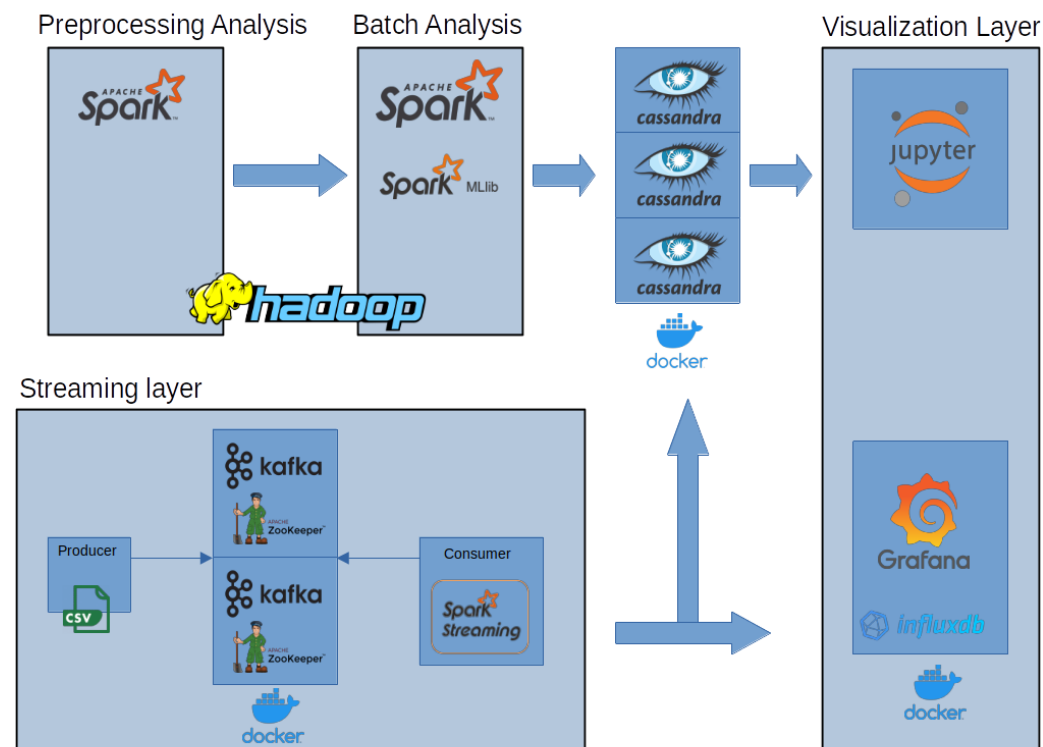
Docker

Portainer

Grafana

InfluxDB

ARCHITETTURA LAMBDA



Containers

Columns Settings

Start Stop Kill Restart Pause Resume Remove Add container

Search...

<input type="checkbox"/>	Name	State Filter	Quick actions	Stack	Image	Created	IP Address	Published Ports	Ownership
<input type="checkbox"/>	cassandra1	running	Logs Info Logs Terminal	cassandra-cluster	bitnami/cassandra:latest	2021-09-26 12:22:23	172.23.0.4	7000:7000 7000:7000 9042:9042 9042:9042	administra
<input type="checkbox"/>	cassandra3	running	Logs Info Logs Terminal	cassandra-cluster	bitnami/cassandra:latest	2021-09-26 12:22:23	172.23.0.2	7002:7000 7002:7000 9044:9042 9044:9042	administra
<input type="checkbox"/>	cassandra2	running	Logs Info Logs Terminal	cassandra-cluster	bitnami/cassandra:latest	2021-09-26 12:22:23	172.23.0.3	7001:7000 7001:7000 9043:9042 9043:9042	administra
<input type="checkbox"/>	portainer	running	Logs Info Logs Terminal	-	portainer/portainer-ce	2021-08-24 11:26:44	172.17.0.2	8000:8000 8000:8000 9001:9000 9001:9000	administra

Items per page 10

DOCKER & PORTAINER



Preprocessing & Batch Layer



Pulizia dei Dati e Merge dei Dataset.

Batch Analysis : analisi effettuate con 4 diverse granularità temporali.

Salvataggio dei report su Cassandra.



Batch Layer: Spark MLlib



Preprocessamento dei dati tramite pipeline: applicazione di label encoder, vettorizzazione e standardizzazione.

Addestramento di un modello per predire voli in ritardo/cancellati.

STREAMING LAYER



Kafka come modulo di data ingestion,
simulando un feed di live-data.

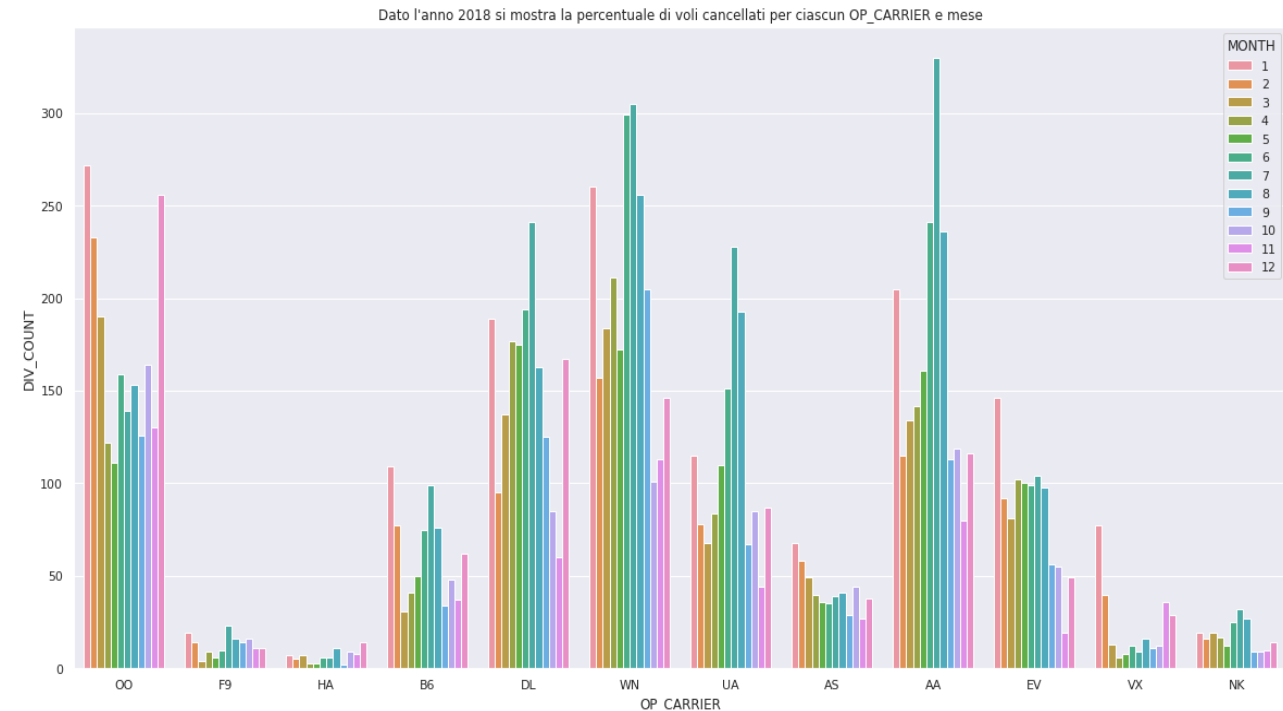
Delay Report salvati su Cassandra.

Preprocessing dei Dati Grezzi salvati
sottoforma di Time-Series su InfluxDB
per successiva visualizzazione.

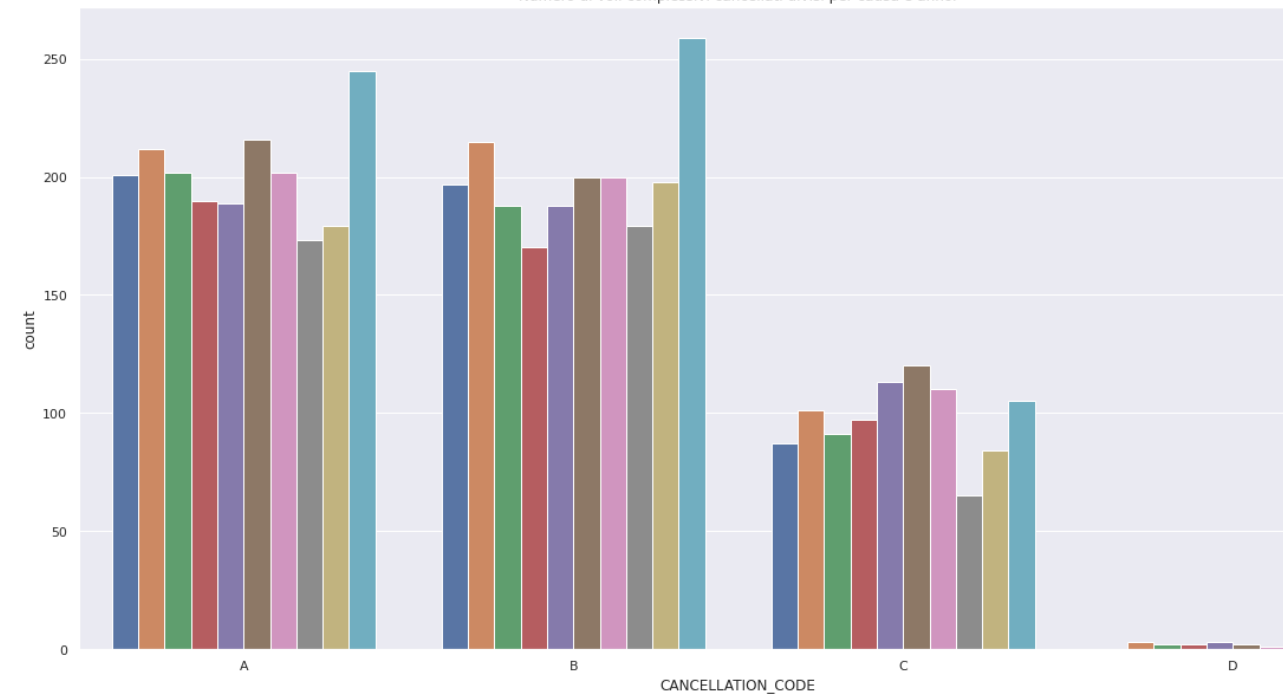
VISUALIZATION LAYER

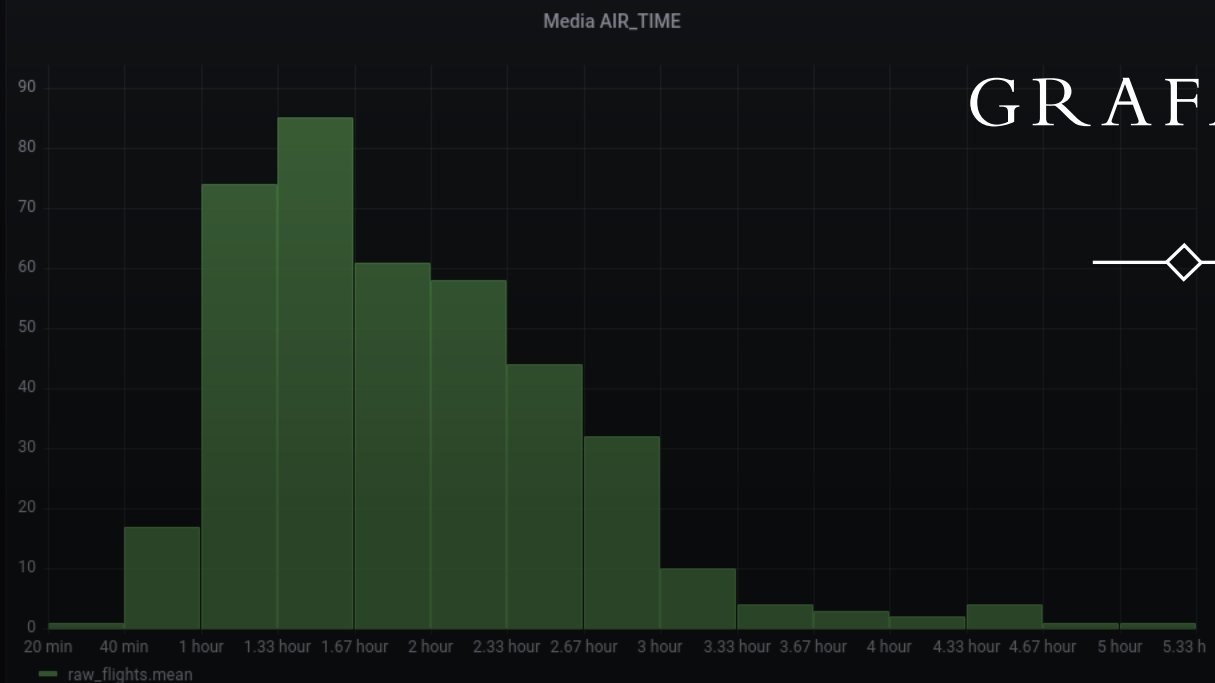
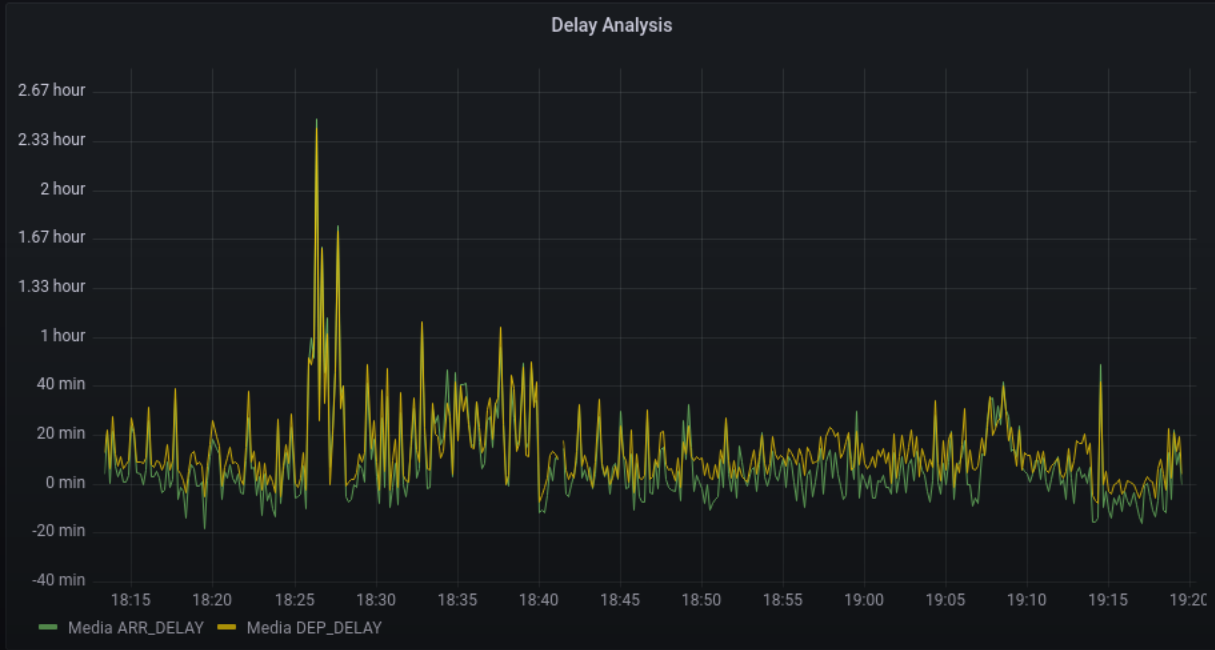


JUPYTER FOR BATCH ANALYSIS VISUALIZATION



Numero di voli complessivi cancellati divisi per causa e anno.





GRAFANA

Flights

Time	raw_flights.AIR_TIME
2021-09-25 18:13:26	38 min
2021-09-25 18:13:26	1.57 hour
2021-09-25 18:13:26	42 min
2021-09-25 18:13:26	31 min
2021-09-25 18:13:26	1.68 hour
2021-09-25 18:13:26	44 min
2021-09-25 18:13:26	1.25 hour
2021-09-25 18:13:26	49 min
2021-09-25 18:13:26	1.15 hour
2021-09-25 18:13:26	3.42 hour
2021-09-25 18:13:26	1.20 hour
2021-09-25 18:13:26	1.25 hour

raw_flights.AIR_TIME

GRAZIE PER
L'ATTENZIONE!

