



Inspirating ideas, connecting people

Daniele Gotti - 1079011 Filippo Bolis - 1079493



...

## Job Pyspark – Load-Data

```
# JOB ADD WATCH NEXT:
## LETTURA DEI DATI DALLA TABELLA "related videos.csv" + JOIN AL DATASET MAIN DEI WATCH-NEXT
# 1. Collego l'URI S3 dei dati che si trovano nel bucket
watch_next_dataset_path = "s3://prog-connectedx-data/related videos.csv"
# 2. Leggo i dati dal csv usando option per gestire la lettura su diverse righe
watch next dataset = spark.read \
    .option("header", "true") \
    .option("quote", "\"") \
    .option("escape", "\"") \
    .csv(watch next dataset path)
# 3. Seleziona le colonne utili all'app dalla tabella dei video collegati
watch next dataset = watch next_dataset.groupBy(col("id").alias("id ref")).agg(
    collect_list(struct(
        col("related id").alias("id"),
        col("title").alias("title"),
        col("presenterDisplayName").alias("presenter")
    )).alias("related videos")
# 4. Svolge il join tra il dataset principale e il nuovo watch next dataset
tedx dataset main = tedx dataset main.join(watch next dataset, tedx dataset main.id == watch
    .drop("id ref")
tedx dataset main.printSchema()
```

I dataset, contenenti informazioni sui video TEDx, vengono uniti, arricchiti con dettagli, immagini e video correlati e infine scritti in un database MongoDB per un facile accesso e analisi

```
slug: "ryan_panchadsaram_anjali_grover_and_david_biello_an_updated_action_pla..."
 speakers: "Rvan Panchadsaram, Aniali Grover and David Biello"
 title: "An updated action plan for solving the climate crisis - and a look at ..."
 url: "https://www.ted.com/talks/ryan_panchadsaram_anjali_grover_and_david_bi..."
 description: "When it comes to climate, what are we doing right and where should we ..."
 duration: "722"
 publishedAt: "2024-04-29T13:28:48Z"
 url image: "https://talkstar-photos.s3.amazonaws.com/uploads/585acd28-d676-449a-a0..."
related_videos : Array (6)

→ 0: Object

      id: "83767"
      title: "An action plan for solving the climate crisis"
      presenter: "John Doerr and Rvan Panchadsaram"
  ▶ 1: Object
  ▶ 2: Object
  ▶ 3: Object
  ▶ 4: Object
  ▶ 5: Object

▼ tags: Array (8)
   0: "climate change"
   1: "environment"
```

Codice completo Click qui



Il codice utilizza AWS Glue con PySpark per caricare dati da un nuovo file CSV su S3 relativo ai video visti dagli utenti. I dati vengono aggregati per l'id dell'utente, creando una lista di video visti per ciascuno

```
1 id_user,id_video
2 0,52688
3 0,5285
4 0,52778
5 1,52688

_id: "2"

video_visti: Array (5)
0: "52688"
1: "528495"
2: "527782"
3: "526878"
```

**4:** "522927"

```
### AGGREGO I DATI PER id_user
user_video_agg = user_video_dataset.groupBy(col("id_user")).agg(
    collect_list(col("id_video")).alias("video_visti")
### SELEZIONO SOLO LE COLONNE id user E video visti
user video agg = user video agg.select(col("id user").alias(" id"), "video visti")
### STAMPA LO SCHEMA DEL DATASET AGGREGATO
user_video_agg.printSchema()
### CONVERTE IN DYNAMICFRAME E SCRIVI IN MONGODB
write_mongo_options = {
    "connectionName": "CONNECTEDX_Connection",
    "database": "connectedx_database",
    "collection": "user_video",
    "ssl": "true",
    "ssl.domain_match": "false"
```

Codice completo Click qui



## Job Pyspark – *User-Data*

```
### PROCESSO I DATI
user dataset = user dataset.withColumn(" id", col("id")).drop("id")
### AGGREGO T DATT VIDEO
user_video_agg = user_video_dataset.groupBy(col("id_user")).agg(
    collect list(col("id video")).alias("video visti")
user_video_agg = user_video_agg.select(col("id_user").alias("_id"), "video_visti")
### "ESPLODO" VIDEO VISTI PER FARE IL JOIN CON I TAG
user_video_exploded = user_video_agg.withColumn("id_video", explode(col("video_visti")))
### JOIN TRA TAGS DATASET CON VIDEO DATA "ESPLOSI"
tags with user = tags dataset.join(user video exploded, tags dataset.id == user video exploded.id video, "inner")
### FACCIO LA CONTA DEI TAGS PER USER
tags_count_per_user = tags_with_user.groupBy(user_video_exploded["_id"], tags_dataset["tag"]).agg(
    count("*").alias("tag count")
### RI-AGGREGO IL TUTTO PER PRENDERE I TAGS E LA LORO CONTA
tags aggregated = tags count per user.groupBy(col(" id")).agg(
    collect list(struct(col("tag"), col("tag count").alias("tag count"))).alias("tags viewed")
### JOIN USER DATA CON VIDEO DATA
user data full = user dataset.join(user video agg, " id", "left")
### JOIN CON AGGREGATED TAGS DATASET
user data full = user data full.join(tags aggregated, " id", "left")
```

#### Codice completo su Click qui

Aggrega i video visti dagli utenti, calcola la frequenza di visualizzazione dei tag e unisce queste informazioni con il dataset degli utenti

```
id: "2"
 name: "Gabriele"
  surname : "Mazzoleni"
 position: "Mariano (BG)"
 coordinateX: "45.640513"
 coordinateY: "9.580923"
 email: "mazzoleni@gmail.com"
 password: "9101"
 url_user_img: "https://static-00.ic
▼ video_visti : Array (5)
    0: "52688"
    1: "528495"
    2: "527782"
    3: "526878"
    4: "522927"
▼ tags_viewed : Array (27)
  ▼ 0: Object
      tag: "math"
```

tag count: 1



Le API di Eventbite che volevamo utilizzare per ottenere i dati degli eventi sono state disabilitate nel 2019

# **Eventbrite API v3**

Disabled in December 12, 2019



Abbiamo quindi optato per un operazione di Scraping del sito



## Per effettuare lo scraping abbiamo:

01

Aperto il sito di Eventbrite

04

Dal sito *Stevesie* rimappato e aggregato i dati

02

Premuto F12 per ispezionare la pagina

05

Copiato i dati in un file Excel eliminando le colonne a noi inutili 03

Scaricato i dati Json inviati degli eventi

06

Convertito il file in csv tramite una macro e caricato su S3



### Abbiamo ottenuto 248 dati:

4	Α	E	3	С	D		E	F	G		Н	1		J	K		L	М
1	name,lan	nguage,	event_ı	url,summ	ary,img_	url,start	_date,t	icket_url	,price,en	d_date	,city,lon	gitude,l	atitud	e,addres	s			
2	COLTIVA	RE L'UM	1ANO S	EMINARI	O RESIDE	NZIALE	en-us,	https://v	vww.even	tbrite.	it/e/colti	ivare-lu	mano-	-seminar	io-reside	nziale	-tickets-	9238999887
3	Bergamo	Outdoo	or Esca	pe Game	: Atalanta	a the He	ro,en-ι	ıs,https:/	//www.ev	entbrit	e.com/e	e/berga	mo-ou	ıtdoor-es	cape-gai	ne-at	alanta-th	ne-hero-ticke
4	Crema - \	Norksh	op Foto	ografia Ri	tratto,it-i	t,https:/	/www.	eventbrit	e.it/e/bigl	ietti-c	rema-wo	orkshop	-fotog	rafia-ritr	atto-7783	30680	5607,wo	rkshop in for
5	SUMMER	PARTN	ER,en-	us,https:	//www.e	ventbrit	e.com/	e/summ	er-partne	r-ticke	ts-9277	703350	47,IL F	PRIMO PA	ARTY DEL	LA CO	NSULEN	IZA ITALIANA
6	Luminous	s Garde	ns - Ale	essandro	Martire,	en-us,ht	tps://w	ww.ever	ntbrite.it/e	e/lumii	nous-gar	dens-a	lessar	ndro-mar	tire-ticke	ts-898	8651429	617,Il pianis
7	INNERHE	RO® - II	Viaggi	o dell'Ero	e per il tu	o succe	sso ne	lla vita,it-	it,https://	www.	eventbrit	te.it/e/r	egistra	azione-in	nerhero-	il-viag	gio-delle	eroe-per-il-tı
8	Fusion 36	60   Cor	so Con	npleto,it-	it,https://	/www.e	entbrit/	te.it/e/bi	glietti-fus	ion-36	0-corso-	-comple	eto-93	8189218	207,Il co	rso af	fronta le	tecniche di r
9	Super sur	nday,en	-us,htt	ps://www	w.eventb	rite.it/e/	super-	sunday-t	ickets-87	84232	16487,E	vento fo	ormati	vo di 5 or	e sul netv	vork,,	2024-09	-15,https://w
10	Ramada	Summe	r Festiv	val 2024,	it-it,https	://www.	eventb	rite.it/e/	biglietti-ra	mada	-summe	r-festiv	al-202	4-91384	6157457	,8 ser	ate sotto	le stelle 8 ba
11	Coffee Ta	ılk: A co	lazione	e con Gea	a Smith, it	it,https-	://www	eventbr.	rite.it/e/bi	glietti-	coffee-t	alk-a-co	olazior	ne-con-g	ea-smith	-9323	5308216	67,Coffee Tal
12	Shakunta	la - Il dr	amma	indiano a	Milano,i	t-it,http:	s://ww	w.eventb	rite.it/e/b	iglietti	-shakun	tala-il-d	Iramm	na-indian	o-a-mila	no-95	6736294	01,Shakunta
13	Il Codice	delle Er	nozion	i & del Co	rpo con i	l dottor l	Bradley	Nelson,	en-us,htt	os://w	ww.even	tbrite.it	t/e/il-c	odice-de	elle-emoz	zioni-c	del-corpo	o-con-il-dott
14	DRINK &	FLOWE	RS,en-	us,https:	//www.e	entbrite/	e.it/e/d	rink-flow	ers-ticke	ts-932	7014942	277,Flor	ral Sips	s Bloomii	ng Momei	nts,ht	tps://img	g.evbuc.com
15	Welcome	Summ	er! Ape	eritivo di ı	networkir	ıg,it-it	,https:/	//www.e	ventbrite.	it/e/bi	glietti-we	elcome-	summ	ner-aperi	tivo-di-ne	etwork	king-932	402048627,
16	Performa	nce Pa	rco La	mbro,en	us,https:	//www.	eventb	rite.com	/e/perfori	mance	-parco-l	ambro-	ticket	s-205824	1876747,	perfo	rmance a	acrobatica p
17	CANDLE	PAINTIN	VG,en-	us,https:	//www.e	ventbrite	e.it/e/c	andle-pa	inting-tic	kets-9	3267134	14097,0	Color Y	our Cano	lle Light Y	our W	orld, http://orld	os://img.evbi
18	POT PAIN	ITING,e	n-us,h	ttps://wv	w.event	orite.it/e	/pot-p	ainting-ti	ckets-932	26655	36727,P	ot Paint	ing: W	here Cre	ativity Blo	oms,	https://ir	mg.evbuc.co
19	Evento Ca	ashflow	Club M	1ilano 21	Settemb	re 2024,	it-it,htt	ps://www	w.eventbr	ite.it/e	/biglietti	i-evento	-cash	flow-clu	b-milano	-21-se	ettembre	-2024-9268
20	TWO DAY	'S,en-u	s,https	://www.e	ventbrite	.com/e	/two-da	ays-ticke	ts-920129	97017	07,Due g	iorni di	forma	zione pe	r la tua cr	escita	profess	ionale,https:
21	Techno B	eats 00	1,en-u	s,https://	www.eve	entbrite.	com/e	techno-	beats-00	1-ticke	ts-9124	103830	17,Te	chno Bea	its is the f	irst EC	OM event	in Bolzano,ł
22	WORSHO	P: Ges	tione c	ontabile	del condo	minio s	enza st	ress,en-	us,https:/	/www	.eventbr	ite.it/e/	worsh	op-gestic	one-cont	abile-	del-cond	lominio-senz
23																		1528397,Dip
24	DRINK & I	BEADIN	G,en-ι	us,https:/	/www.ev	entbrite	.it/e/dr	ink-beac	ling-ticke	ts-932	6738515	97,Wh	ere Be	ads Beco	ome Mast	terpie	ces,http:	s://img.evbu



## Job Pyspark – *Event-Data*

```
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'])
### LEGGO IL FILE CSV PER CREARE IL DATASET DI INPUT
events_dataset = spark.read \
    .option("header", "true") \
    .option("quote", "\"") \
   .option("escape", "\"") \
    .csv(events_dataset_path)
### STAMPO LO SCHEMA DEL DATASET
events dataset.printSchema()
### CONVERTE IN DYNAMICFRAME E SCRIVI IN MONGODB
write mongo options = {
    "connectionName": "CONNECTEDX Connection",
    "database": "connectedx_database",
   "collection": "event_data",
    "ssl": "true",
    "ssl.domain match": "false"
### CONVERTE IN DYNAMICERAME
events dynamic frame = DynamicFrame.fromDF(events dataset, glueContext, "nested")
### SCRIVO IN MONGODB UTILIZZANDO LE OPZIONI SPECIFICATE
glueContext.write_dynamic_frame.from_options(events_dynamic_frame, connection_type="mongodt
```

Dopo aver letto e verificato lo schema del dataset, lo converte in un DynamicFrame e lo scrive in un database MongoDB

```
_id: ObjectId('668d56f5bc57800dff417938')
name: "COLTIVARE L'UMANO SEMINARIO RESIDENZIALE"
language: "en-us"
event_url: "https://www.eventbrite.it/e/coltivare-
summary: "Percorsi di autodeterminazione per un appi
img_url: "https://img.evbuc.com/https%3A%2F%2Fcdn.or
start_date: "2024-08-17"
ticket_url: "https://www.eventbrite.com/checkout-es
price: "139.93"
end_date: "2024-08-20"
city: "Barzio"
longitude: "9.4531160999999999"
latitude: "45.9424051"
address: "Valsassina 23816 Barzio"
```

Codice completo Click qui





Dati statici degli eventi

In mancanza dell'API di Eventbrite siamo stati costretti a svolgere uno scraping, tuttavia i dati ottenuti sono utili solo nel breve periodo

**Password esposte** 

Le password degli utenti sono pubbliche nel database MongoDB – S3

Doppio DB

In caso di aggiornamento dei dati su S3, è necessario riavviare ogni volta il job di AWS Glue, il che può causare ritardi.

For more info: GitHub | Trello







Dati statici degli eventi

Se in futuro Eventbrite offrirà API per l'accesso ai dati degli eventi, sarà possibile implementarle. In alternativa, si potrebbe considerare il cambio del provider dei dati.

**Password esposte** 

Si potrebbe rendere l'accesso più sicuro, ad esempio con Cognito





