
CSCI-E88C Final Project

Charles Lariviere

2022-12-08

Contents

1	[50 points] Project Proposal and Summary	3
1.1	Project Goal and Problem Statement	3
1.2	Data Sources	3
1.3	Expected Results	4
1.4	Application Overview and Technologies Used	5
1.5	Bonus Options	5
2	[50 points] Milestone 1: System Design and Architecture	6
2.1	System Diagram	6
2.2	Description of each stage of the processing pipeline	6
3	[100 points] Milestone 2: Implementation and Results	7
3.1	a. Code and Configuration Snippets	7
3.1.1	Configuration	7
3.1.2	EventGenerator	8
3.1.3	Events	8
3.1.4	Message / Messageable	9
3.1.5	SparkStreamingAggregator	10
3.1.6	EventTypeAggregator	11
3.1.7	ZeroResultSearchesAggregator	12
3.1.8	Logstash	13
3.2	Screenshots of run commands and processing results	14

1 [50 points] Project Proposal and Summary

1.1 Project Goal and Problem Statement

Creative Market is an online marketplace that allows creatives to buy and sell digital design assets, such as fonts, templates, photos, and web themes. Users can browse a catalog of millions of digital assets and purchase licenses that allow them the use of the asset, with prices ranging from a few dollars to a few thousand dollars.

To find assets, users can browse the various curated collections and pages, or use the search engine to find the specific type of asset they are looking for. Given the large size of the catalog, content delivery and personalization is crucial to provide a great user experience and return highly-relevance assets that are likely to meet the user's needs and result in a purchase.

To achieve this goal, the Machine Learning team at Creative Market (which I am a part of) develops multiple systems and models that power all content delivery on the website, with the goal of optimizing relevance and revenue. While there exists evaluation frameworks across some of these initiatives, visibility on the general content delivery performance is limited.

The goal of this project is to create a stream processing application that will visualize multiple performance metrics, based on various user interaction events with the results returned by the various locations on the website. This will provide a robust understanding of the overall performance of the Machine Learning team's initiatives across the various content delivery systems.

1.2 Data Sources

For every search query performed by a user, a set of events are recorded to track the user's interactions with the results. For instance, if a user performs a search for "business card mockups" and receives a grid of 30 results, an event will be recorded for each of the result returned to the user.

Furthermore, as the user browses through the set of results and interacts with them – views, clicks, adds to cart, and purchases – another event is saved for each, recording the product identifier, the user identifier, as well as a unique identifier representing the search query.

These events are sent to a stream in Amazon Kinesis, and then stored in Amazon S3 in batches of compressed JSON files partitioned by type and date using Amazon Kinesis Data Firehose.

Here is an example payload for an event:

```
{
  "type": "ProductView",
  "search_id": "256c8e20-3fb0-4409-b56e-307fcf99790c",
  "user_id": "f73f9219-146a-4a21-8123-c945181d7cb3",
  "product_id": "aaa621de-40fe-446a-83bf-e5fcbee7fd3a",
  "recorded_at": "2022-11-10T13:43:10.35462",
  "version": "1.3.71"
}
```

Figure 1: ProductView Event

The types of events recorded are:

- Search: Recorded for each search request made by a user, with the payload including the search parameters, filters, and terms.
- Result: Recorded for each product returned in a search result set.
- ProductView: Recorded every time a product included in a search result is viewed (i.e. enters the user's viewport in their browser window).
- ProductClick: Recorded every time a product included in a search result is clicked on.
- ProductAddtoCart: Recorded every time a product included in a search result is added to a user's cart.
- ProductPurchase: Recorded every time a product included in a search result is purchased.

1.3 Expected Results

The output of this stream processing application should be a set of time-series metrics assessing the performance of content delivery across the website, grouped by various dimensions.

Typical evaluation metrics for information retrieval will be used, such as click-through rate, Success Rate, zero-result rate, and Normalized Discounted Cumulative Gain (nDCG). Here are the definitions for each of the metrics that will be computed:

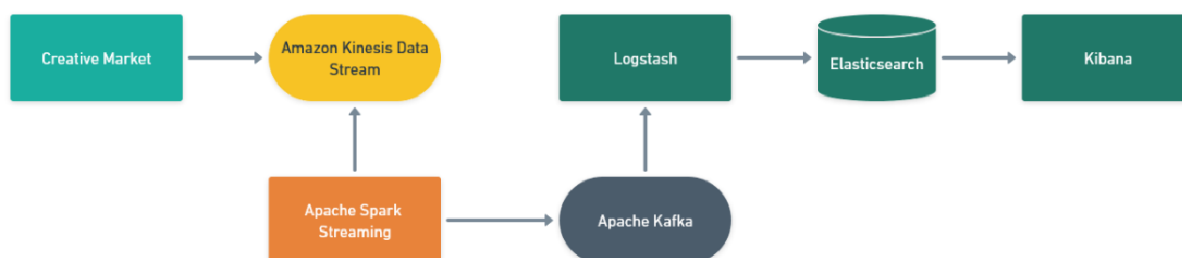
- Click-Through Rate: Ratio of results viewed by a user that were then clicked on.
- Success Rate: Ratio of searches that resulted in a purchase.

- Zero-Result Rate: Ratio of searches that did not return any results.
- nDCG: Score ranging from 0 to 1 that represents how the relevance of actual search results compared to the ideal ranking. Here, click-through rate is used as a proxy for relevance.

The dimensions made available for analysis will include the date, type of page (i.e. Search Engine, Category Page, “You may also like” module), user country, and application version number (i.e. it indicates which ML model and ranking function was used).

These metrics will be stored in Elasticsearch and made available for consumption in Kibana.

1.4 Application Overview and Technologies Used



The streaming application will consume from the stream of events in Amazon Kinesis and leverage Apache Spark Streaming to execute the transformation.

The transformation will be windowed and consist of joining events together on the original result event, to produce a record for each result along with descriptive attributes used to compute the metrics (i.e. isViewed, isClicked, isPurchased).

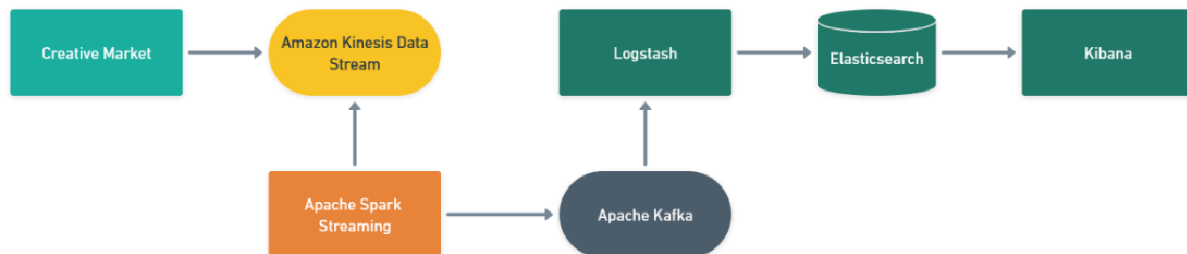
The transformed records will then be written to a Kafka stream which is used to push these to ELK using the Kafka logs integration. Records will be saved in Elasticsearch which will be used to aggregate the metrics and visualize them in Kibana.

1.5 Bonus Options

I will attempt to cover all three bonus options, by writing unit tests on the transformation executed in Spark, executing the streaming application on AWS by consuming from a Kinesis stream, and creating visualizations in Kibana,

2 [50 points] Milestone 1: System Design and Architecture

2.1 System Diagram



2.2 Description of each stage of the processing pipeline

- **Collection:** Events are streamed to Amazon Kinesis by an external system and made available in an Amazon Kinesis Data Stream. Events are made available for consumption within a few milliseconds and are cached for a period of 7 days.
- **Ingestion:** Events will be consumed from the Kinesis data stream using Spark Streaming's Kinesis integration. It allows for obtaining a Dstream from a Kinesis stream by simply providing the necessary configuration.
- **Preparation:** Very limited preparation is needed since the events are formatted as well-formatted JSON objects.
- **Computation:** A multi-stage transformation process is required for this application. Given a stream of ProductSearch, ProductResult, ProductView, ProductClick, ProductAddtoCart, and ProductPurchase, we want to return a set of calculated metrics per window. Computing metrics like Click-Through-Rate, Success Rate, and nDCG will require events to be joined together on the ProductID key, to generate enrich the original ProductResult object with isViewed, isClicked, and isPurchased attributes. Once joined and enriched, metrics such as Click-Through-Rate and Success Rate can be computed per window.
- **Presentation:** The computed metrics will be made available for consumption and visualization by using the ELK stack – Elasticsearch, Logstash, and Kibana. Once computed and logged into the Kafka topic, Logstash's Kafka integration will automatically consume the records and store them into Elasticsearch. These will then be visualized through charts and dashboards using Kibana.

3 [100 points] Milestone 2: Implementation and Results

3.1 a. Code and Configuration Snippets

3.1.1 Configuration

Configuration for the Scala portion of this project is provided through a `.conf` file and Scala case classes.

```
1  org.cscie88c {
2    kafka {
3      bootstrap-server = "localhost:9092"
4    }
5    spark {
6      master-url = "local[*]"
7    }
8    event-generator {
9      name = "product-event-generator"
10     topic = "product-events"
11     searches-per-minute = 6
12     view-rate = 0.7
13     click-rate = 0.3
14     purchase-rate = 0.1
15   }
16   count-aggregator {
17     name = "count-aggregator"
18     topic = "product-event-counts"
19   }
20   zero-result-searches-aggregator {
21     name = "zero-result-searches-aggregator"
22     topic = "product-zero-result-searches"
23   }
24 }
```

```
1  final case class KafkaConfig(bootstrapServer: String)
2  final case class SparkConfig(masterUrl: String)
3  final case class EventGeneratorConfig(
4    name: String,
5    topic: String,
6    viewRate: Float,
7    clickRate: Float,
8    purchaseRate: Float,
9    searchesPerMinute: Int
10 )
11 final case class CountAggregatorConfig(name: String, topic: String)
12 final case class ZeroResultSearchesAggregatorConfig(name: String, topic: String)
```

These configuration case classes are then loaded using `pureconfig` in each object that requires them. For instance;

```
1 val KAFKA_CONFIG_PATH = "org.cscie88c.kafka"
2 implicit val kafkaConf: KafkaConfig = ConfigSource.default.at(KAFKA_CONFIG_PATH)
   ).loadOrThrow[KafkaConfig]
```

3.1.2 EventGenerator

While, in practice, the events are generated by an external service and made available in an Amazon Kinesis Data Stream, I decided to instead generate synthetic data for the purpose of this project. This synthetic data is written to a Kafka topic that is then consumed by the Spark Structured Streaming applications.

```
1 object EventGenerator extends App {
2   ...
3   val producer = new KafkaProducer[String, String](properties)
4
5   // generate an infinite stream of ProductSearch instances
6   val searchStream: Stream[ProductSearch] = Stream.continually(ProductSearch.
      create())
7
8   // for each item, wait, generate records, and send to a Kafka topic
9   searchStream
10    .map(sleep(_, 60000 / eventConf.searchesPerMinute))
11    .map(getRecords)
12    .foreach(x => x.map(log).foreach(producer.send))
13
14   def getRecords(productSearch: ProductSearch): List[ProducerRecord[String,
      String]] = {
15     // each ProductSearch event generates multiple ProductSearch,
16     // ProductView, ProductClick, and ProductPurchase events; each
17     // with its own configured probability
18     val results = createResults(productSearch)
19     val views = results.flatMap(_.createView(eventConf.viewRate))
20     val clicks = views.flatMap(_.createClick(eventConf.clickRate))
21     val purchase = clicks.flatMap(_.createPurchase(eventConf.purchaseRate))
22
23     // concatenate all events
24     val events = List(productSearch) ::: views ::: clicks ::: purchase
25
26     // generate an instance of ProductRecord for each event/message
27     events.map(_.toMessage.toRecord(eventConf.topic))
28   }
29   ...
30 }
```

3.1.3 Events

Each of the events generated by this class are structured as Scala case classes with a companion object used as an instance factory. For example, the `ProductResult` event is implemented as such:


```
1 case class ProductResult(id: String, searchId: String, productId: String, rank:
  Int, recordedAt: String) extends Messageable {
2   override def key: String = this.id
3   override def timestamp: String = this.recordedAt
4   override def toJson: Json = this.asJson
5
6   def createView(probability: Float): Option[ProductView] = {
7     // create a ProductView with a certain probability
8     if (Random.nextFloat() < probability) {
9       Option(ProductView.create(searchId = this.searchId, productId = this.
        productId))
10    } else {
11      None
12    }
13  }
14 }
15
16 object ProductResult extends Helper {
17   def create(searchId: String, rank: Int): ProductResult = {
18     // create a ProductResult instance with random attributes
19     ProductResult.apply(
20       id = randomId,
21       searchId = searchId,
22       productId = randomId,
23       rank = rank,
24       recordedAt = currentDate
25     )
26   }
27 }
```

This object takes advantage of a few powerful Scala features, such as Option, Trait, and Implicits.

3.1.4 Message / Messageable

To standardize the structure of messages sent to Kafka, and to avoid duplicating logic in each of the multiple `Product*` case classes, a `Message` class and `Messageable` trait were created.

```
1 case class Message(key: String, eventType: String, payload: String, recordedAt:
  String) {
2   def toRecord(topic: String): ProducerRecord[String, String] = {
3     new ProducerRecord[String, String](topic, key, this.asJson.noSpaces)
4   }
5 }
```

```
1 trait Messageable {
2   def key: String
3   def eventType: String = {
4     this.getClass.getSimpleName
5   }
6   def toJson: Json
7   def timestamp: String
8   def toMessage: Message = {
```

```
9      // create an instance of Message from the current class
10     Message(
11         key = this.key,
12         eventType = this.eventType,
13         payload = this.toJson.noSpaces,
14         recordedAt = this.timestamp
15     )
16 }
17 }
```

3.1.5 SparkStreamingAggregator

With the synthetic data available in a Kafka topic, it is now possible to run our Spark Structured Streaming applications. Given that the scope of this project includes multiple types of aggregations, a `SparkStreamingAggregator` Trait was created to facilitate the creation of multiple Spark Structured Streaming applications that rely on Kafka by reducing code duplication.

```
1  trait SparkStreamingAggregator extends App {
2      val SPARK_CONFIG_PATH: String = "org.cscie88c.spark"
3      val KAFKA_CONFIG_PATH: String = "org.cscie88c.kafka"
4
5      // load config
6      implicit val sparkConfig: SparkConfig = loadConfig[SparkConfig](
7          SPARK_CONFIG_PATH)
8      implicit val kafkaConfig: KafkaConfig = loadConfig[KafkaConfig](
9          KAFKA_CONFIG_PATH)
10
11     // create spark session
12     implicit val spark: SparkSession = SparkSession
13         .builder()
14         .master(sparkConfig.masterUrl)
15         .getOrCreate()
16
17     // import spark implicits
18     import spark.implicits._
19     spark.sparkContext.setLogLevel("ERROR")
20     ...
21 }
```

While this Trait initializes the required config and `SparkSession`, it also provides a few powerful utility methods that can be used by downstream classes.

For example, it provides a `getEvents` method that returns a `DataFrame` from the Kafka topic that contains our synthetic data, and unpacks the JSON payload using the schema of the `Message` class defined above (i.e. `Encoders.product[Message].schema`)

```
1  def getEvents(topic: String)(implicit spark: SparkSession, kafkaConf:
2      KafkaConfig): DataFrame = {
3      val df = spark
4          .readStream
5          .format("kafka")
```

```
5      .option("kafka.bootstrap.servers", kafkaConf.bootstrapServer)
6      .option("subscribe", topic)
7      .option("startingOffsets", "earliest") // From starting
8      .load()
9
10     // decode key/value into string
11     df.select(
12         col("key").cast("string"),
13         col("value").cast("string")
14     )
15     // extract JSON payload using Message schema
16     .select(
17         from_json(
18             col("value"),
19             Encoders.product[Message].schema
20         ).as("parsed")
21     )
22     // unpack parsed JSON into columns
23     .select("parsed.*")
24     .withColumn("recordedAt", col("recordedAt").cast("timestamp"))
25 }
```

3.1.6 EventTypeAggregator

Given the logic contained in the Trait defined above, our Spark Structured Streaming application that is responsible for the actual business logic is relatively slim. This application counts the number of events by 1-minute window and type of event.

```
1  object EventTypeAggregator extends SparkStreamingAggregator {
2      val COUNT_AGGREGATOR_CONFIG: String = "org.cscie88c.count-aggregator"
3      val EVENT_GENERATOR_CONFIG: String = "org.cscie88c.event-generator"
4      implicit val config: CountAggregatorConfig = loadConfig[CountAggregatorConfig]
5      implicit val generatorConfig: EventGeneratorConfig = loadConfig[EventGeneratorConfig]
6
7      // get events from Kafka topic
8      val events = getEvents(generatorConfig.topic)
9
10     // count by window and eventType
11     val counts = events
12         .withWatermark("recordedAt", "1 minute")
13         .groupBy(
14             window(col("recordedAt"), "1 minute"),
15             col("eventType")
16         )
17         .count()
18
19     // format as Kafka messages
20     val df = toKafkaDataFrame(counts, col("window").cast("string"))
21
22     // write to Kafka topic
```

```
23 writeKafka(df, config.topic)
24 }
```

3.1.7 ZeroResultSearchesAggregator

The following application filters and parses events of type `ProductSearch` from the topic that contains all types of events, and then counts the number of `ProductSearch` where `results=0`.

```
1 object ZeroResultSearchesAggregator extends SparkStreamingAggregator {
2   val ZERO_RESULT_AGGREGATOR_CONFIG: String = "org.cscie88c.zero-result-
   searches-aggregator"
3   val EVENT_GENERATOR_CONFIG: String = "org.cscie88c.event-generator"
4   implicit val config: ZeroResultSearchesAggregatorConfig = loadConfig[
   ZeroResultSearchesAggregatorConfig](ZERO_RESULT_AGGREGATOR_CONFIG)
5   implicit val generatorConfig: EventGeneratorConfig = loadConfig[
   EventGeneratorConfig](EVENT_GENERATOR_CONFIG)
6
7   // get events from Kafka topic
8   val events = getEvents(generatorConfig.topic)
9   // parse JSON into columns using schema from case class
10  val searches = getEventsOfType[ProductSearch](events)
11
12  val zeroResultSearches = searches
13    .filter(col("results").equalTo(0))
14
15  // count number of rows where 'results' = 0
16  val counts = zeroResultSearches
17    .withWatermark("recordedAt", "1 minute")
18    .groupBy(
19      window(col("recordedAt"), "1 minute"),
20    )
21    .count()
22
23  // write to Kafka topic
24  val df = toKafkaDataFrame(counts, col("window").cast("string"))
25  writeKafka(df, config.topic)
26 }
```

This application takes advantage of the `getEventsOfType` polymorphic method defined in the `SparkStreamingAggregator` trait. This method allows for “automatically” parsing and unpacking a JSON payload into DataFrame columns with the correct types.

```
1 def getEventsOfType[T <: Product : ClassTag : TypeTag](events: DataFrame):
   DataFrame = {
2   events
3   // filter rows where 'type' is equal to the name of the provided class
4   .filter(col("eventType").equalTo(lit(classTag[T].runtimeClass.
   getSimpleName)))
5   // extract JSON payload using T schema
6   .select(
7     from_json(
```

```
8         col("payload"),
9         Encoders.product[T].schema
10        ).as("parsed")
11    )
12    // unpack parsed JSON into columns
13    .select("parsed.*")
14    .withColumn("recordedAt", col("recordedAt").cast("timestamp"))
15 }
```

3.1.8 Logstash

Finally, since the Spark streaming applications above write their output to a Kafka stream, we can take advantage of the ELK stack's integration with Kafka to automatically import messages into Elasticsearch using Logstash.

The logstash pipeline that imports messages output by `EventTypeAggregator` is as follows:

```
1 input {
2     kafka {
3         codec => json
4         bootstrap_servers => "kafka:29092"
5         topics => ["product-event-counts"]
6     }
7 }
8 output {
9     elasticsearch {
10         hosts => ["elasticsearch:9200"]
11         user => "elastic"
12         password => "changeme"
13         index => "product-event-counts"
14     }
15 }
```

3.2 Screenshots of run commands and processing results

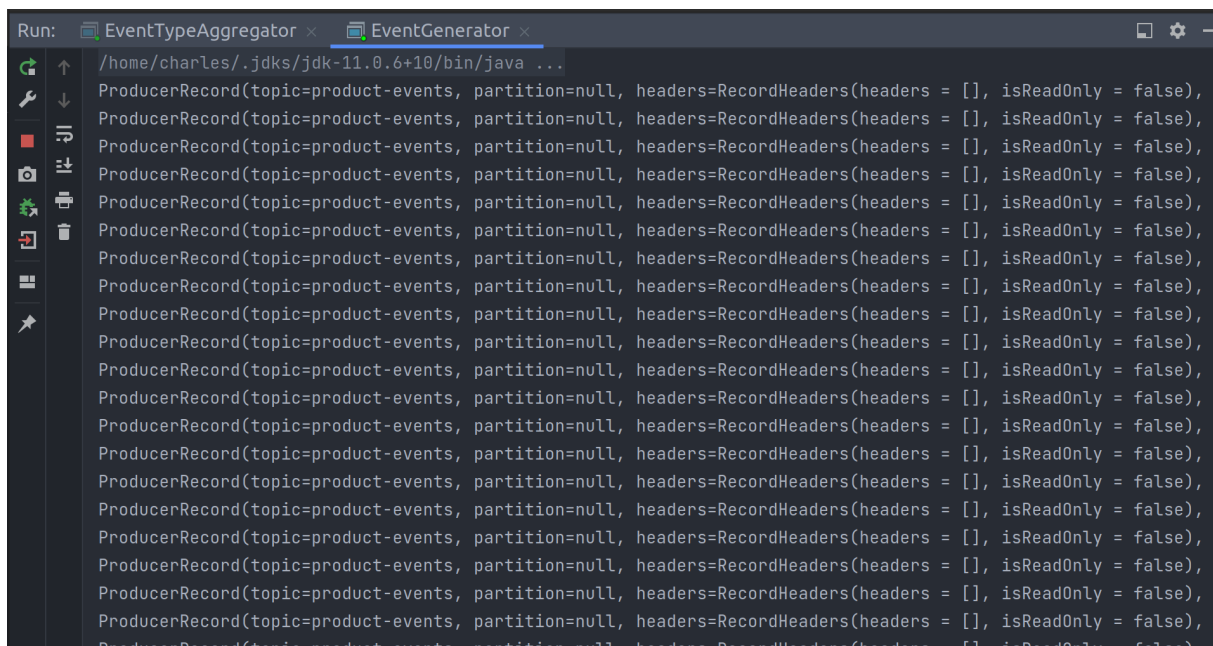


Figure 2: ProducerRecord generated by EventGenerator

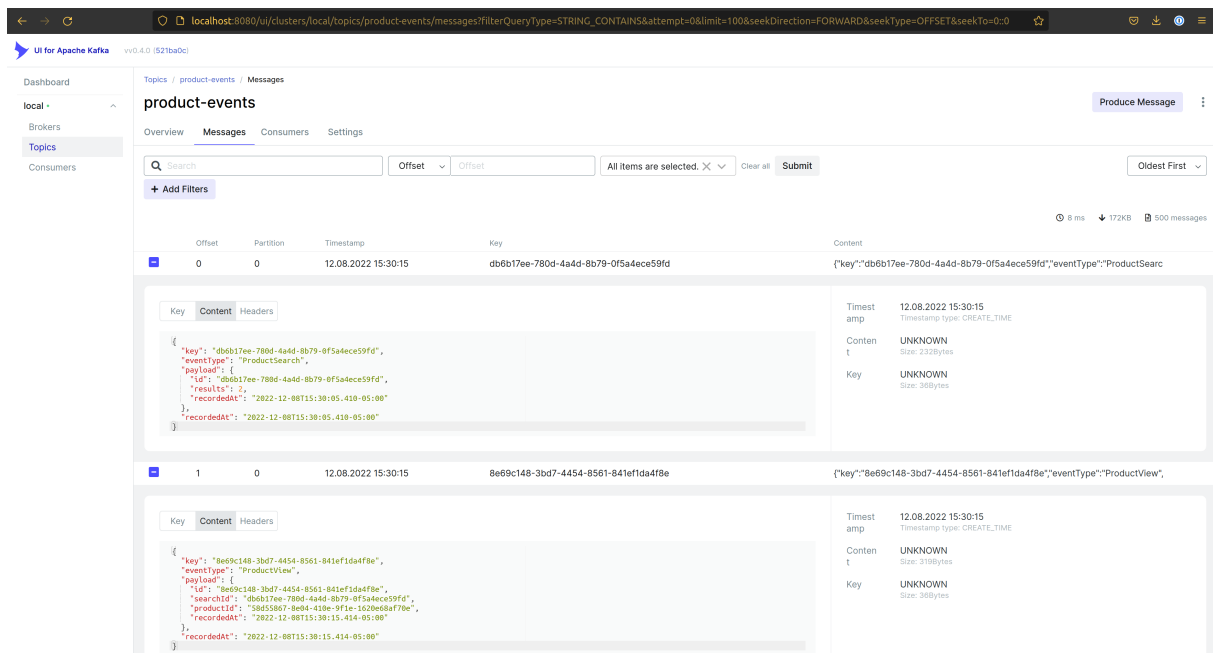


Figure 3: Output of EventGenerator in Kafka

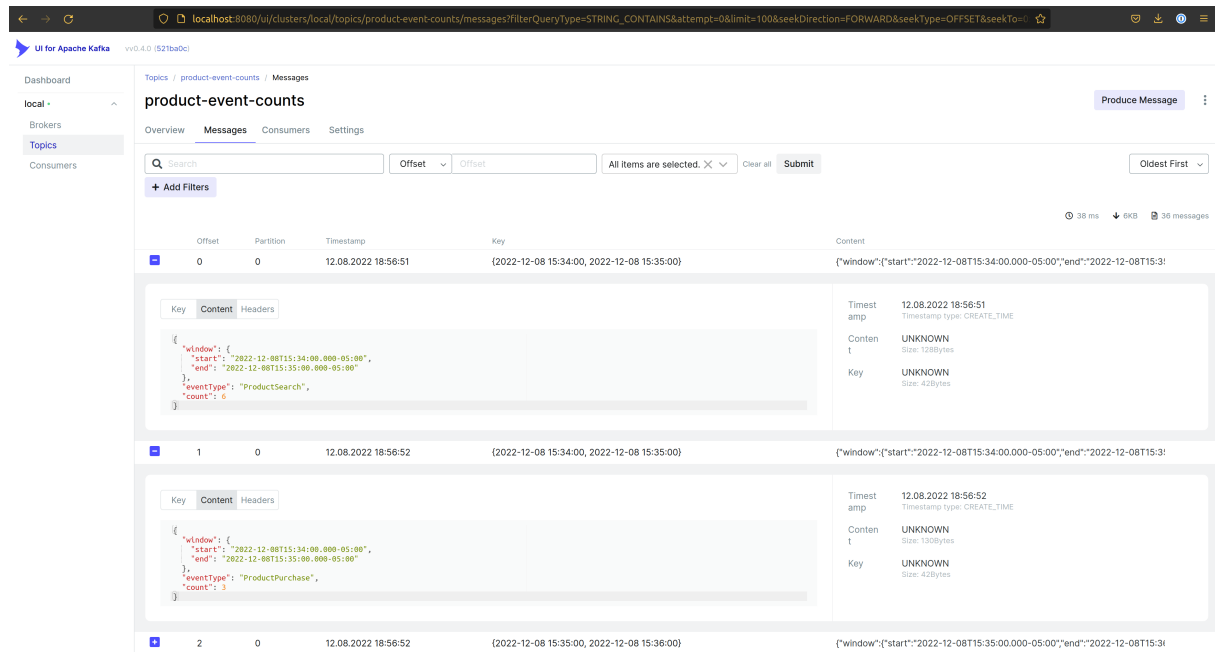


Figure 4: Output of EventTypeAggregator in Kafka

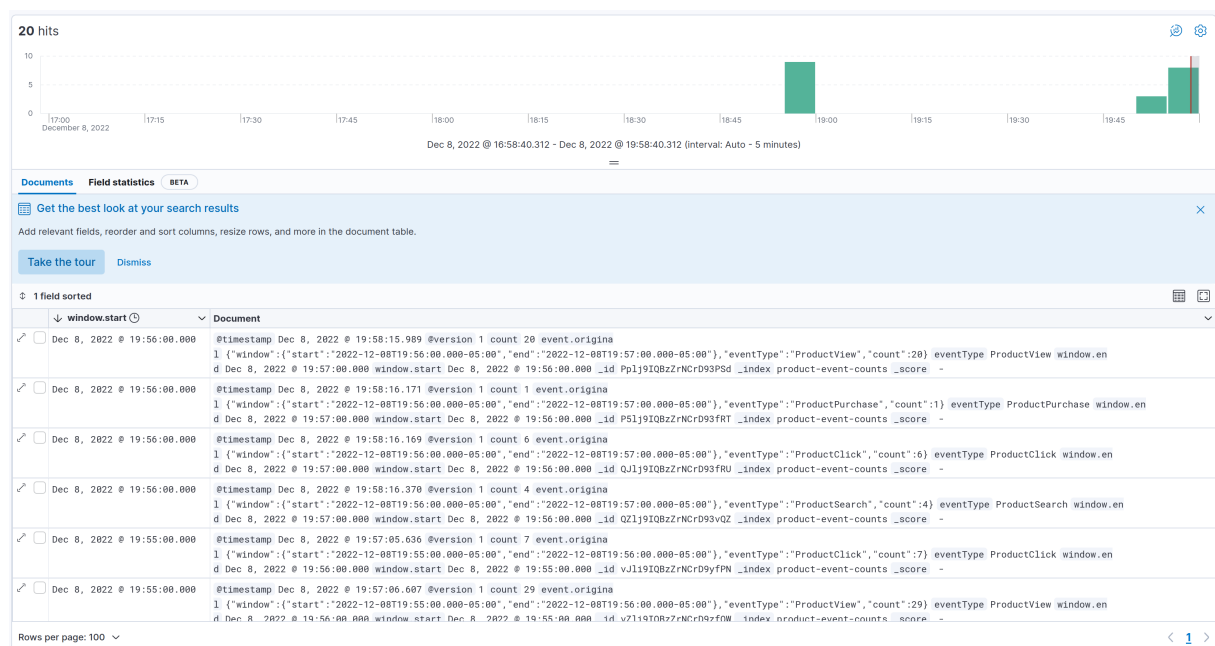


Figure 5: Output from EventTypeAggregator in Kibana