# Project Report

# A Dataflow Journey
# from PubSub to BigQuery and Data Studio

## Google Cloud Pipeline
### Data Engineering 2

## Elnaz Dehkharghani
Student Number: 11015404
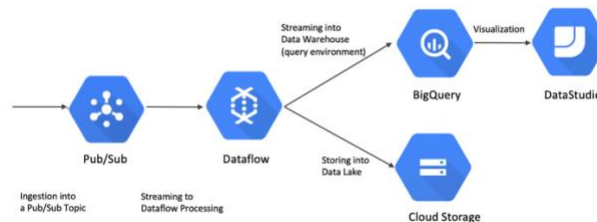M.Sc. Big Data and Business Analytics

## Prof. Dr. Frank Schulz

SRH Heidelberg University
16th Jun 2021

**Introduction.**

In Data Engineering one we made our data Pipeline in a manual way. These kinds of pipelines are great for doing some quick prototypes and exploratory analysis or collecting data in a fast way from somewhere. In contrast, this time in Data Engineering two, the goal is to look at a larger scenario and production qualities. We would have a scalable pipeline and as a result, will more consider the big data. In this schema we guaranty the on timing and availability of the system while people are using it. The whole process will be done in the cloud. So, the main goal of this project is to move from prototype level to production level.
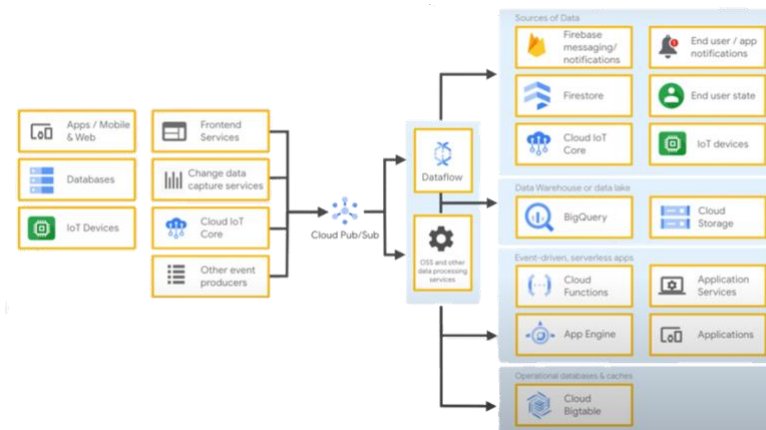
**The schema of the desired pipeline:**



Before deep-diving into data flowing and storytelling, we will have a brief description of all the pipeline components.

1. **Pub/Sub:** This is our messaging component. If you are ingesting large amount of data for analysis, or looking to simplify the development of event-driven microservices then you need cloud pub/sub. It is similar to Kafka that has topics. Pub/Sub helps to build robust, scalable systems of applications by integrating them asynchronously. Cloud pub/sub is a fully managed, real-time messaging service that allows you to send and receive messages between independent applications. Once an event is successfully published as a message, it becomes the job of the service to ensure that all the systems that need to react to this event get it. This architecture is more scalable and reduces dependencies. Pub/Sub can store messages for seven days, so in case of some temporary failures in some parts, they can receive it yet. As the name pub/sub indicates, cloud pub/sub supports a publisher-subscriber model. A publisher application creates and sends messages to a topic which is named resource and the messages are stored until acknowledged by all the subscribers. To receive these messages a subscriber application creates a subscription to a topic. The subscriber receives a message by either cloud pub/sub pushing them to the subscribers chosen endpoint or by subscriber pulling them from the service. When the message is acknowledged by the subscriber it will remove from the subscription backlog and not delivered again. Communicating can be one-to-many (fan-out), many-to-one (fan-in), and many-to-many.

   - **Publishers:** can be any application that can make https request to Google APIs.com.
   - **Subscribers:** Pull subscribers can also be any application although push subscribers must be a webhook endpoint that can accept post requests.

Example of incoming data sources to Cloud Pub/Sub and its targets

2. **Dataflow:** In this pipeline, dataflow is actually Apache Beam that is an open-source tool. It can do many tasks such as select, transfer, transformation, and enriching data with additional information in the streaming way "Record base". We have to consider streaming here does not mean that our data is live because they are actually stored in a pub/sub topics. So, dataflow sends data message by message or record by record. This part actually acts as the consumer. Based on Apache Beam, this Google Cloud service is used for data processing both in batch or streaming mode using the same code, providing horizontal scalability to calibrate the resources needed. The service provides many templates ready to use. Among them, there is one, that follows the pipeline schema that I am going to elaborate on later: reading messages from Pub/Sub and writing them to BigQuery.

3. **BigQuery:** Storing and querying massive datasets can be time-consuming and expensive without the right hardware and infrastructure. BigQuery is an enterprise data warehouse that solves this problem by enabling super-fast SQL queries using the processing power of Google's infrastructure.
It is possible to control access to both the project and data based on business needs, such as giving others the ability to view or query your data. BigQuery is accessible by using the cloud console, by using the bq command-line tool, or by making calls to the BigQuery REST API using a variety of client libraries such as Java, .NET, or Python. There are also a variety of third-party tools that you can use to interact with BigQuery, such as visualizing the data.

4. **Cloud Storage:** For the longtime storing data we need this component as a data lake. We can also use some alternative cloud data lakes.

5. **Data Studio:** For the visualization aspect, we can use Google Data Studio. This is one environment that we can easily create our dashboards and visualize our data.

**Advantages of using this Pipeline:**
- Very fast implementation and setup
- Scalable
- Availability
- Robust
- Working with larger data

**Description of Data Source and Dataset: London Crime 2008-2016**

Now that we have a brief overview of our data pipeline, it is time to make familiarize ourselves with our data and making some user stories to know how to get useful insights from our data.

For this project, I am going to use London crime data between the years 2008 to 2016, which is in CSV file format and has around 13.5 million rows and 7 columns. The dataset is downloaded from the Kaggle website. Data is already cleaned and ready for doing analysis.
This data counts the number of crimes at two different geographic levels of London (Lower Layer Super Output Area "LSOA" and Boroughs) by year, according to crime type. This is telling us, if the table naming convention is sensible, that this data is going to be grouped primarily by the area where the crime/s happened.

**Columns are as following:**

- **lsoa_code:** This represents a policing area
- **borough:** The London borough for which the statistic is related
- **major_category:** The major crime category
- **minor_category:** The minor crime category
- **value:** The count of the crime for that particular borough, in that particular month
- **year**: The year of the summary statistic
- **month**: The month of the summary statistic

**User story telling:**

1. How many boroughs are there in London?
2. What is the rate of crimes in each borough?
3. How many LSOA codes are there per borough?
4. How many crimes have happened per year in London from 2008 to 2016?
5. What is the overall volume of the different crimes in London?
6. How are the crimes in Croydon over the whole period?
7. What is the order of major crime categories for Croydon in 2012 that most of the crimes have happened in that year?
8. What type of crimes is growing fastest?
9. What are the different types of crimes and their quantity per year?
10. What are the minor categories in Theft and Handling?
11. What is the change in the number of crime incidents from 2011 to 2016?
12. What are the top 3 crimes per borough?

### Steps for Making the Data Pipeline in Google Cloud Console

**Step One:** Creating a GCP Account

It is really simple just by having one Gmail account everybody can have a Google Cloud account as well. There is a good option for every individual when joining Google Cloud to have 300 USD free credit for 90 days to play around, gaining experience, and then upgrade their accounts into payable ones.

**Interacting with GCP account:**
1. **Web Interface**
2. **Command Line Interface:** For using the CLI should install Google Cloud SDK base on your machine hardware name then initialize it and continue using by gcloud utility, with connecting it to your google account. For the BigQuery you can also use another utility called bq. Here I have successfully installed the GCP SDK and initialized it.



**Step Two:** Make a Project in GCP

Creating the project that we are going to deal with it in the whole process here named as DataEngineering. It is very important to save the Project ID somewhere since we will need it in further steps.



**Step Three:** Enabling the Required APIs

For this project, we have to enable the Dataflow API and Cloud Pub/Sub API.



Since Google Cloud Platform has lots of tools, it is better to pin your required components at the beginning to have better access of them.

**Step Four:** Create a Topic

In this step, it is time to create our topic for sending and receiving the messages. We do this by referring to the Pub/Sub section. When we create our topic, below of it a long topic name is generated that is always accessible in the Pub/Sub section. We will need this also in further steps to make our dataflow. For the sake of this project, we do not need to make subscriptions but if we had some subscriber app, we need to define them in the subscription menu as many as that we have subscribers.



**Step Five:** Create a Dataset and Tables

We do this by referring to the BigQuery section and add our dataset and tables there.



When we are going to make our table, it is very important to have a good understanding of our table schema otherwise during the sending messages will encounter problems. So, to do a better job we can use this command in the terminal and see our CSV file records there:

$ cat london_crime_by_lsoa.csv

We can copy some of the records and paste them into the JSON convertor website to see that how our messages will be sent in the JSON format.

JSON

[
  {
    "lsoa_code": "E01001116",
    "borough": "Croydon",
    "major_category": "Burglary",
    "minor_category": "Burglary in Other Buildings",
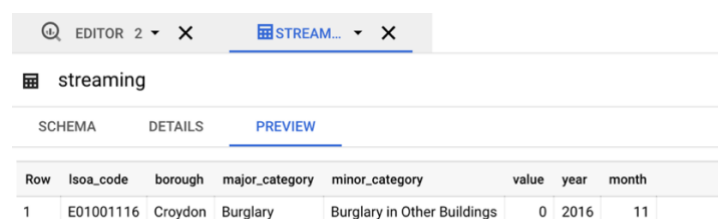    "value": 0,
    "year": 2016,
    "month": 11
  },

First four of the columns should be in string type and the last three should be in integer. Now that we have enough confidence about our data types it is time to make table schema. Since in this step I am going to send one-sample message to my table I chose an empty table with its required schema.

As soon as that we make our table, one table ID is generated for us which is very important to use in the dataflow configuration step.

**Step Six:** Create a Bucket in Google Cloud Storage

The only thing that we need is just choosing a bucket name nothing else. I chose it the same as my project name. There are some more options for handling your bucket but this will affect the monthly fees.

**Step Seven:** Create a Dataflow

Dataflow can be assumed as our consumer. We set this with the provided template as Pub/Sub topic to BigQuery.



It is exactly here that we need to add our topic name, BigQuery table ID and our bucket address location.

When we run the job, we will get this dataflow chart that is ready to listen to our publisher and store our data into BigQuery.



To examine if everything works well, we send one message manually in Pub/Sub "Publish Message" section. The message should be in JSON format.



Now we refer to the BigQuery section to see if our message is there, and yes, it is there.

As far as everything works well, it is time to either write a python script as our publisher to send data from there or store our data file locally in Google Cloud Storage and transferring them by another job template "Text Cloud Storage to PubSub Topic".

**Making Publisher Application in Python:** Building a One-to-One Pub/Sub System

In the case of making a publisher app, we need to set up a service account and Cloud IAM. This is important to our application to use it for authentication and set up Cloud IAM permissions. For this purpose, we navigate to the IAM & Admin section then select the Service accounts. We choose our service account name and click create. After that, should set the roles that we need for our apps either Pub/Sub Publisher or Pub/Sub Subscriber or maybe both. For this example, the service account needs both publishing and subscribing permissions.



By having the service account, can go inside its link to create a key. This key will be used by the client library to access the Cloud Pub/Sub API. We select JSON and create it.



From this step to further it is very important to have our Google Cloud SDK installed. In continue for sending and receiving sample messages through pub.py and sub.py we should clone the repository of the Google Cloud in our specified folder and then create a virtual environment and activate it.

To examine this python publisher and subscriber I have made a new topic named "hello_topic" and made a subscriber for it as "sub_one" that pulls the messages from the topic.



**Further Steps:**
1. Download or clone the required files in local folder:
   $ git clone https://github.com/googleapis/python-pubsub.git

2. Make a virtual environment and activate it.

3. Download the required files:
   $ pip install --upgrade google-cloud-pubsub

4. Set the project that are going to use:
   $ gcloud config set project dataengineering-315913

5. Watching all topics and subscriptions inside the project:
   $ gcloud pubsub topics list
   $ gcloud pubsub subscriptions list

6. Doing the Google Credentials:
   $ export GOOGLE_APPLICATION_CREDENTIALS=~/Desktop/Data_Engineering/keys/key_test.json

7. Setting the value:
   $ export PROJECT=`gcloud config get-value project`

8. Opening two new terminal that have the pub.py and sub.py inside it.



Python Programs for Publishing and Receiving Messages

9. Running the subscriber program:
   $ python3 sub.py $PROJECT sub_one



Subscriber is ready to listen to messages.

10. Finally in another terminal run the publisher program:
    $ python pub.py $PROJECT hello_topic



Cloud pub/Sub receive the message ID that indicate the subscriber has received the message and do not need to resend it.
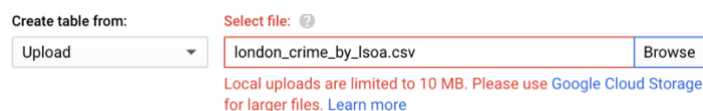


Now that we have a short experience in using python programs for sending and receiving messages through Pub/Sub it is time to send our big data from Pub/Sub topic to BigQuery. We can do this with a little change in our python program to read our CSV file record by record or in JSONL format, and by keep running one job template just like before as "PubSub Topic to BigQuery" write our real messages in BigQuery.

**Step Eight:** Accessing to Our London Crime Data for Doing the BigQuery in Four Different Attempts Through GCP Console

**8.1 First Attempt - Uploading the Dataset from Local System to BigQuery:**

In the previous steps, we have seen that I have already created my "londoncrime dataset" so the only thing that we need is to make a new table. Since my data file with 13.5 million records is 932,8 MB an error message is pop up that indicates to use google cloud storage for your large files. Here I decided to make my file shorter and try this option from my local system and after that continue with my origin file from google cloud storage.



Therefore, I selected only 5000 rows from the total records and continue my job:

This time schema and file format is automatically detected. The data is inserted successfully and ready for doing queries. As I have decided to make my queries with larger amount of data in the next attempt, I will upload my origin data to Google Cloud and from there insert them to BigQuery.

## 8.2 Second Attempt - Uploading the CSV to Google Cloud Storage and ingesting in BigQuery:

For this purpose, we refer to the created bucket with the name of dataengineering-315913 and upload our CSV file with 13.5 rows there.



Now the CSV file is actually there and ready for ingesting by BigQuery. Further, navigate to the BigQuery section and make a new table. This time instead of the upload option we select the Google Cloud Storage to upload our original file with 13.5 million rows.



The file format and table schema are recognized automatically and we do not need to do anything. We have our data in two tables, one from the local uploading and the other from Google Cloud Storage:

## 8.3 Attempt Three - Transferring Data to BigQuery by The Help of one Dataflow and GCS:

For this option we have to do a little programming in JavaScript to make the UDF function. Since I am using the free version of google cloud, in this step will use my smaller file with only 5000 rows.

- **8.3.1 Making a small UDF program in JavaScript:**



- **8.3.2 Making the JSON schema of the CSV file:**



- **8.3.3 Uploading the UDF, JSON Schema and my small CSV file without header to the GCS.**



- **8.3.4 Creating a target empty table with defining its schema:**

- **8.3.5 Create a job template from dataflow section:** Text Files on Cloud Storage to BigQuery (Batch)

We select "Text Files on Cloud Storage to BigQuery (Batch)" template and then address the location of our UDF function, name of the function, schema of our table in JSON format, targeted BigQuery Table, our supposed CSV file without headers and two temporary directories then run the job. In this process everything goes well except inserting data into BigQuery table.



## 8.4 Attempt Four – Transferring Data to BigQuery with the help of Two Dataflow Using the Pub/Sub:

At this point with the help of Pub/Sub I can insert my data to BigQuery table successfully.

Steps for this purpose are as followings:

- **8.4.1 Making a new empty table with its defined schema:**

- **8.4.2 Using the JSONL file from the last step:**

The previous step (8.3-Attempt three) has generated a beautiful JSON line file in our temporary folder to use it in this step.



We refer to the provided URL and copy the generated JSONL file into our local system with a .json extension. It is important that for this step and even the previous step omit the header of the files. After that we upload this JSON file to the Google Cloud Storage.

- **8.4.3 Create two different dataflow job templates:**
  **8.4.3.1 First Dataflow: Pub/Sub Topic to BigQuery**

Here we fill the fields with our topic which we have created it at the beginning of this document.



Then address the new table and finally make a new temporary folder for it and run the job.

We have to keep this job running to can send our messages through topic.



**8.4.3.2 Second Dataflow: Text File on Cloud Storage to Pub/Sub Topic**
We have our listener, so it is time for sending our messages and we do it by making a new job template from dataflow. This time we choose "Text File on Cloud Storage to Pub/Sub" template to send all our 5000 messages through it. It is important that from the GCS select our JSONL without header otherwise we will encounter with errors in the BigQuery that cannot serialize our messages in the proper way.

As soon as that we run this job 200 messages from the whole are visible in the BigQuery table. For finalizing the message sending I had to keep this job running about 2 hours to see all 5000 rows in my BigQuery table.

## Step Nine: BigQuery
**BigQuery and answering the predefined user stories:**

Now that we could successfully store our data in four different ways into BigQuery it is time to do different queries and visualize them to have a good insight of our data. We do these queries on the data that we inserted from our GCS to BigQuery which it has 13.490.604 rows in total. We can immediately see that there are more than one lsoa_code per borough telling us that this represents a smaller more confined geographical area within a borough; this may seem straightforward and obvious, but the relationship between the variables here might help us pose some questions that we can later answer using the information available to us.

```
-- 1. How many boroughs are there in London ?
SELECT DISTINCT borough
FROM `dataengineering-315913.londoncrime.fromGCE_13M`;


-- 2. What is the rate of crimes in each borough?
SELECT borough, SUM(value) AS number_of_crimes
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
GROUP BY borough
ORDER BY number_of_crimes DESC;


-- 3. How many LSOA codes are there per borough?
SELECT borough, count(DISTINCT lsoa_code) as n_codes
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
GROUP BY borough
ORDER BY count(DISTINCT lsoa_code) DESC;


-- 4. How many crimes have happend per year in London from 2008 to 2016?
SELECT year, SUM(value) AS `total_crime`
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
WHERE year BETWEEN 2008 AND 2016
GROUP BY year
ORDER BY total_crime DESC;


-- 5. What is the overall volume of the different crimes in London?
```

```sql
SELECT year, major_category, SUM(value) AS `total_crime`
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
GROUP BY year, major_category
ORDER BY year, total_crime DESC;
```

-- 6. How are the crimes in Croydon over the whole time period?

```sql
SELECT year, SUM(value) AS `total_crime`
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
WHERE borough = 'Croydon'
GROUP BY year
ORDER BY total_crime DESC;
```

-- 7. What is the order of major crime categories for Croydon in 2012 that most of the crimes has happened in that year?

```sql
SELECT major_category, sum(value) AS `total_crime`
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
WHERE borough = 'Croydon' AND year = 2012
GROUP BY major_category
ORDER BY total_crime DESC;
```

-- 8. What type of crimes are growing fastest? (by finding the total crime counts for the major categories across all the years in the data set)

```sql
SELECT year, major_category, SUM(value) AS `total_crime`
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
GROUP BY major_category, year
ORDER BY major_category, year;
```

-- 9. What are the different types of crimes and their quantity per year?

```sql
SELECT year, major_category, SUM(value) AS `total_crime`
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
WHERE year BETWEEN 2008 AND 2016
GROUP BY major_category, year
ORDER BY major_category, year;
```

-- 10. What are the minor categories in Theft and Handling?

```sql
SELECT major_category, minor_category, year,SUM(value) AS total_possesions
FROM `dataengineering-315913.londoncrime.fromGCE_13M`
WHERE major_category='Theft and Handling'
GROUP BY 1,2,3
ORDER BY 3 DESC;
```

-- 11. What is the change in the number of crime incidents from 2011 to 2016?

```sql
SELECT
 borough,
 no_crimes_2011,
 no_crimes_2016,
 no_crimes_2016 - no_crimes_2011 AS change,
```

```
    ROUND(((no_crimes_2016 - no_crimes_2011) / no_crimes_2016) * 100, 2) AS perc_change
FROM (
  SELECT
    borough,
    SUM(IF(year=2011, value, NULL)) no_crimes_2011,
    SUM(IF(year=2016, value, NULL)) no_crimes_2016
  FROM
    `dataengineering-315913.londoncrime.fromGCE_13M`
  GROUP BY
    borough )
ORDER BY
  perc_change ASC;


-- 12. What are the top 3 crimes per borough?
SELECT
  borough,
  major_category,
  rank_per_borough,
  no_of_incidents
FROM (
  SELECT
    borough,
    major_category,
    RANK() OVER(PARTITION BY borough ORDER BY SUM(value) DESC) AS rank_per_borough,
    SUM(value) AS no_of_incidents
  FROM
    `dataengineering-315913.londoncrime.fromGCE_13M`
  GROUP BY
    borough,
    major_category )
WHERE
  rank_per_borough <= 3
ORDER BY
  borough,
  rank_per_borough;
```

## Step Ten: Data Visualization with Data Studio:

Visualization of the data is possible in two different ways. First making query and do the visualization at the same time in the data studio in the Explore Data section, the other way is to initially refer to the Google Data Studio make a blank report and from there connecting your BigQuery and project then do different kind of queries in a single report as your dashboard. There is also one option to share your work with others and managed it with different kind of permissions. In the following I did my visualization in data studio report and answered to all above queries with numbering them.

**Link to the Dashboard:** https://datastudio.google.com/reporting/ac4051cd-f515-4209-992d-f19b3e741916

## London Crime Data Set
### 2008 - 2016 Geographic Crime Case

Created by: Elnaz Dehkharghani | Data Source Kaggle.com | 16th June 2021

**(1) Total Number of Boroughs in London**

Record Count
33

**(2) Top 5 Boroughs by Crime Rate**



number_of_crimes

| Borough | |
|---|---|
| Westminster | 455,028 |
| Lambeth | 292,178 |
| Southwark | 278,809 |
| Camden | 275,147 |
| Newham | 262,024 |

▷ **Westminster** is the borough with the highest crime rate; the **City of London** is the borough with the lowest crime rate.

**(3) Top 5 Boroughs with Most LSOA Code**



n_codes

| Croydon | Barnet | Bromley | Ealing | Enfield |
|---|---|---|---|---|
| 220 | 211 | 197 | 196 | 183 |

▷ **Croydon** is the borough with the highest number of LSOA codes; the **City of London** is the borough with the lowest.

**(4) Total Number Committed Crimes each Year in London**



total_crime

| 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|---|---|---|
| 738,641 | 717,214 | 715,324 | 724,915 | 737,329 | 686,407 | 680,183 | 711,624 | 736,121 |

▷ It is obvious between 2008 to 2016, **2008** is the year with the highest amount of crimes and in contrast to the other years, **2014** is the more controlled one.

**(5) Volume and Type of the Crimes**



- Theft and Handling
- Violence Against the Person
- Burglary
- Criminal Damage
- Drugs
- Robbery
- Other Notifiable Offences
- Fraud or Forgery
- Sexual Offences

▷ In London, **Theft and Handling** are occurring most of the time, after that **violence against the person** and **Burglary** are two categories that are happening mostly respectively.

**(6) Crimes in Croydon Borough Over the Whole Time Period (2008 - 2016)**



total_crime

▷ Most of the crimes has happened in **2012**.

**(7) Major Catogory Crimes in Croydon in 2012**



total_crime   total_crime

▷ In **2012**, there is not any report for both fraud or sexual offenses but instead, in the overall view, **theft and handling** are the most reported cases.

## (8) Analysis of Growth of Different Type of Crimes Over the Whole Time Period



**Theft and Handling** — **Violence Against the Person** — **Burglary** — **Criminal Damage** — **Drugs** — **Robbery** — **Other Notifiable Offences** — **Fraud or Forgery** — **Sexual Offences**

From 2008 to 2016 **other notifiable offenses** and **violence against the person** are those which are growing fastest over time. With this speed, violence against the person is going to reach theft and handling rates soon in the future.

All the other crimes appear to be in reasonable decline; although note the scaling for the y-axis is not truly comparable. We could calculate the percentage change for each of the major categories, which would give us a better indication as to the true problem of crimes that aren't being tackled.

## (9) Different Types of Crimes and Their Quantity per Year



**Theft and Handling** — **Violence Against the Person** — **Burglary** — **Criminal Damage** — **Drugs** — **Robbery** — **Other Notifiable Offences** — **Fraud or Forgery** — **Sexual Offences**

## (10) Minor Categories of Theft and Handling



- **Other theft** is the most committed one and **handling stolen goods** is the least occurred one in this category

## (11) What is the Percent Change in Crime Across Boroughs?
### (From 2011 to 2016)



| Borough | perc_change |
|---|---|
| City of London | 61.24 |
| Hackney | 24.31 |
| Greenwich | 18.38 |
| Tower Hamlets | 15.52 |
| Haringey | 14.7 |
| Bexley | 11.8 |
| Wandsworth | 9.69 |
| Havering | 8.32 |
| Islington | 7.8 |
| Newham | 7.02 |

## (12) Top 3 Crimes Per Borough



**Burglary** — **Violence Against the Person** — **Theft and Handling** — **Criminal Damage** — **Drugs**

**Conclusion:**

Google's Pub/Sub platform is great for handling large amounts of data and decoupling the various components of our architecture. While this project was not a real-life use case, but could have a very good experience and intertwined different Google Cloud Platform services. With a pipeline like this, you can easily migrate your data from applications to serverless, managed services on the cloud. Once the data is in BigQuery, you can do additional analytics and generate visualizations in Google Data Studio.

<div align="right">The End</div>

**Resources:**

[1] Getting Started With Authentication:
https://cloud.google.com/docs/authentication/getting-started
[2] QuickStart: Using Client Libraries:
https://cloud.google.com/pubsub/docs/quickstart-client-libraries
[3] Building a One-to-Many Pub/Sub System:
 https://cloud.google.com/pubsub/docs/building-pubsub-messaging-system
[4] What is BigQuery?
https://cloud.google.com/bigquery/docs/introduction?hl=en_US
[5] Installing Google Cloud SDK:
https://cloud.google.com/sdk/docs/install
[6] Pushing Messages to the Topic:
https://cloud.google.com/pubsub/docs/publisher
[7] Google Cloud Platform Python Samples:
https://github.com/GoogleCloudPlatform/python-docs-samples
[8] Error Handling for Google Package:
https://stackoverflow.com/questions/47011713/importerror-cannot-import-name-pubsub-v1
[9] A Dataflow Journey: From PubSub to BigQuery:
https://medium.com/codex/a-dataflow-journey-from-pubsub-to-bigquery-68eb3270c93
[10] Cloud PubSub Python Sample:
https://github.com/GoogleCloudPlatform/cloud-pubsub-samples-python/blob/master/gce-cmdline-publisher/traffic_pubsub_generator.py
[11] Setting Up a GCP Pub/Sub Integration With Python
http://www.theappliedarchitect.com/setting-up-gcp-pub-sub-integration-with-python/
[12] CSV Streamer with PubSub and BigQuery
https://github.com/thundercomb/csv-pubsub-bq-streamer
[13] Writing data from PubSub to BigQuery via Apache Beam
http://abitdeployed.com/2020/06/08/writing-data-from-pubsub-to-bigquery-via-apache-beam/