

# Học phần Quản trị Dữ liệu lớn: Bài thực hành Apache Kafka (2)

*Phạm Tiến Lâm, Đặng Văn Báo*

## Cài đặt thư viện:

1. Cài đặt **java** trên máy tính (điều kiện tiên quyết)

Link download: <https://www.oracle.com/java/technologies/downloads>

2. Cài đặt **Apache Kafka**

Link download: <https://kafka.apache.org/downloads>

Hướng dẫn cài đặt và sử dụng:

<https://www.youtube.com/watch?v=BwYFuhVhshI&t=358s>

3. Cài đặt **Jupyter notebook** hoặc **Anaconda**

4. Cài đặt **Kafka-Python**

```
pip install kafka-python
```

```
conda install -c conda-forge kafka-python
```

## Chữa bài Lab6:

1. Khởi động máy chủ **Zookeeper**.

```
.\bin\windows\zookeeper-server-start.bat .\config\zookeeper.properties
```

2. Khởi động máy chủ **Kafka**.

```
.\bin\windows\kafka-server-start.bat .\config\server.properties
```

3. Tạo **Topic** mới trong môi trường **Kafka**.

```
kafka-topics.bat --create --bootstrap-server localhost:9092 --replication-factor 1  
--partition 1 --topic "topic_name"
```

4. Mở một ‘**producer console**’ cho phép gửi các message tới **Topic** trong Kafka.

```
kafka-console-producer.bat --broker-list localhost:9092 --topic "topic_name"
```

5. Mở một “**consumer console**” và cho phép đọc tin nhắn từ **Topic**.

```
kafka-console-consumer.bat --topic "topic_name" --bootstrap-server  
localhost:9092 --from-beginning
```

6. Kiểm tra các **Topic** đã khởi tạo.

```
kafka-topics.bat --list --bootstrap-server localhost:9092
```

## Activity 2: Bắt đầu với Python Apache Kafka

### Tạo một jupyter notebook với tên “Producer”

- Khai báo thư viện sử dụng:

```
1 from kafka import KafkaProducer
2 from time import sleep
3 import json
```

- Sử dụng thư viện KafkaProducer để tạo một đối tượng producer, được kết nối đến một Kafka cluster thông qua các máy chủ

```
1 producer = KafkaProducer(bootstrap_servers= ['localhost:9092'], api_version = (0,10,1))
```

- Gửi 1 message tới topic “abcd” (“abcd” là Topic đã khởi tạo trước đó)

```
1 # topic name is abcd
2 producer.send('abcd', b'1')
```

- (\*)Gửi một số nguyên lên Kafka với Topic ‘abcd’, sau khi chuyển đổi giá trị số này sang dạng bytes.

```
1 number_to_send = 20
2
3 # Chuyển đổi số nguyên thành bytes literal
4 number_bytes = str(number_to_send).encode('utf-8')
5
6 # Gửi thông điệp đến Kafka
7 producer.send('abcd', number_bytes)
```

- (\*\*)Đọc một file dữ liệu “u.data” và gửi thông tin các dòng dữ liệu lên topic “abcd”

```
1 import csv
2
3 # Đọc file CSV và gửi từng dòng đến Kafka topic
4 with open('u.data', newline='') as csvfile:
5     reader = csv.reader(csvfile)
6     for row in reader:
7         column_value = row[0]
8         # Gửi từng cột của mỗi dòng đến Kafka topic
9         for column_value in row:
10             producer.send('abcd', value=column_value.encode('utf-8'))
11             print(f"Sent value: {column_value}")
12
```

## Tạo một jupyter notebook với tên “Consumer”

- Khai báo thư viện sử dụng

```
1 from kafka import KafkaConsumer
2 from time import sleep
3 import json
4 from json import loads
```

- Khởi tạo một đối tượng KafkaConsumer để kết nối đến một chủ đề Kafka ('abcd' trong trường hợp này) thông qua các máy chủ Kafka được chỉ định

```
1 consumer = KafkaConsumer('abcd',
2                             bootstrap_servers= ['localhost:9092'],
3                             api_version = (0,10))
```

- Đọc tất cả message nhận được từ Producer gửi tới Topic

```
1 for mess in consumer:
2     print(mess.value)
```

- Từ phần (\*) trong Activity 1: viết 1 đoạn mã lệnh trong Consumer tính bình phương số đó và in ra màn hình kết quả.
- Từ phần (\*\*) trong Activity 1: viết 1 đoạn mã trong Consumer in ra dữ liệu nhận được từ Producer.

## Bài tập:

Sử dụng thư viện **requests** và **BeautifulSoup** lấy thông tin từ trang web  
“<https://vnexpress.net/apple-ra-mat-iphone-14-4508267.html>”

### Yêu cầu:

1. Khởi tạo 1 Producer: tìm ra nội dung thẻ <title> và truyền nội dung thẻ vào topic đã khởi tạo.
2. Khởi tạo một Consumer: in ra nội dung thẻ <title>
3. Khởi tạo 1 Producer: tìm ra nội dung tất cả các thẻ <a> và truyền nội dung thẻ vào topic đã khởi tạo.
4. Khởi tạo một Consumer: in ra nội dung thẻ <a>