

Politechnika Warszawska

WYDZIAŁ ELEKTRONIKI
I TECHNIK INFORMACYJNYCH



Instytut Informatyki

Studia Podyplomowe
Big Data - przetwarzanie i analiza dużych zbiorów danych

PRACA KOŃCOWA

Kamil Gontarz

Zaprojektowanie i wykonanie systemu do
składowania, przetwarzania oraz analizy danych o
aktywności użytkowników strony internetowej.

Opiekun pracy
mgr Patryk Pilarski

Warszawa, 2022

Spis treści

1. Cel biznesowy	3
2. Dostępne dane	4
2.2 Struktura	4
2.3 Statystyki	4
2.2 Generator	4
3. Rozwiązanie koncepcyjne problemu biznesowego	5
4. Model danych	6
5. Opis modułów	7
5.1 Ogólny schemat rozwiązania	7
5.2 Preprocessing	7
5.3 Processing	7
5.4 Storage	7
5.5 Wizualizacja	7
6. Napotkane problemy oraz ich rozwiązanie	8
7. Wnioski	9
8. Bibliografia	10

1. Cel biznesowy

- zwiększenie ilości i zaangażowania użytkowników strony

W wyniku pracy należy: - zaproponować co można zmienić aby zrealizować cel główny - przygotować testowe środowisko: – przechowujące dane – przetwarzające dane – prezentujące dane - środowisko na pierwszym etapie może być testowe ale użyte technologie powinny być łatwo skalowalne, powszechnie używane w środowisku, z rozwiniętą społecznością, w miarę możliwości open source - modułowość - integracja - zastępowalność - jeżeli w pierwszym kroku analiza nie przyniesie efektów to ważne jest aby gotowa była koncepcja systemu przechowującego dane który może posłużyć do dalszych analiz teraz i w przyszłości - ip jako dana osobowa - anonimizacja przy zachowaniu wiedzy na temat ip

2. Dostępne dane

2.2 Struktura

2.3 Statystyki

2.2 Generator

3. Rozwiązanie koncepcyjne problemu biznesowego

- odstony
- sesje
- czas
- podział na segmenty techniczne (os, deviceType, Browser)
- próba znalezienia odpowiedzi na pytanie jak zwiększyć liczbę użytkowników oraz ich zaangażowanie

4. Model danych

- tabela eventów
- tabela sesji

5. Opis modułów

5.1 Ogólny schemat rozwiązania

5.2 Preprocessing

5.3 Processing

5.4 Storage

5.5 Wizualizacja

6. Napotkane problemy oraz ich rozwiązanie

7. Wnioski

8. Bibliografia

[https://altinity.com/blog/2020/5/21/clickhouse-kafka-engine-tutorial – virtual columns](https://altinity.com/blog/2020/5/21/clickhouse-kafka-engine-tutorial-virtual-columns) <https://hub.docker.com/r/clickhouse/clickhouse-docker/> - clickhouse docker <https://github.com/enqueue/metabase-clickhouse-driver> – metabase clickhouse community driver <https://clickhouse.com/docs/en/integrations/kafka/kafka-table-engine/> <https://www.metabase.com/docs/latest/operations-guide/running-metabase-on-docker.html> <https://www.metabase.com/docs/latest/administration-guide/01-managing-databases.html#database-sync-and-analysis> https://vincent.doba.fr/posts/20211004_spark_data_description_language <https://towardsdatascience.com/spark-3-2-session-windowing-feature-for-streaming-data-e404d92e267> <https://pypi.org/project/user-agents/> <https://databricks-prod-cloudfront.cloud.databricks.com/public/4863dbfdce0008a2fedfa335468c0841bb7ad231f2c49cfda2bf1611ee66cc3> <https://stackoverflow.com/questions/70991571/stream-data-from-one-kafka-topic-to-another-using-pyspark> <https://stackoverflow.com/questions/2013124/regex-matching-up-to-the-first-occurrence-of-a-character> <https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html#windows-operations-on-event-time> <https://github.com/ykursadkaya/pyspark-Docker/blob/master/Dockerfile> <https://stackoverflow.com/questions/37132559/add-jar-files-to-a-spark-job-spark-submit> <https://stackoverflow.com/question/do-i-write-to-kafka-using-pyspark> <https://faker.readthedocs.io/en/master/> <https://github.com/joke2k/faker> https://docs.docker.com/develop/develop-images/dockerfile_best-practices/ - multiline env <https://support.google.com/analytics/answer/2731565?hl=en#overview&zippy=%2Cin-this-article-definition-of-visit>