

# Melhorias no Processo e Estratégias Avançadas

## 1. Melhoria no Código Atual

Problemas que podemos enfrentar:

- Extração sequencial dos dados da API pode ser lenta.
- Transformação é feita em memória e pode não escalar bem com grandes volumes de dados.
- Dados transformados são carregados diretamente em um arquivo CSV, o que não é ideal para análises contínuas e armazenamento histórico.

Melhorias sugeridas:

1. **Paralelização das Requisições à API:**
  - Utilizar `concurrent.futures` para realizar requisições em paralelo, reduzindo o tempo de extração dos dados.
2. **Uso de Bibliotecas Eficientes para Manipulação de Dados:**
  - Utilizar bibliotecas como `pandas` para manipulação eficiente de grandes volumes de dados.
3. **Armazenamento Incremental:**
  - Modificar o código para permitir armazenamento incremental, carregando os dados em partes para um data lake ao invés de armazenar tudo em memória.

## 2. Carregamento dos Dados em um Data Lake

Estratégia:

- Após a extração e transformação, os dados podem ser carregados em um banco de dados centralizado (data lake). Este banco de dados pode ser um sistema distribuído como Amazon S3, Google Cloud Storage ou um banco de dados SQL de alta performance como Amazon Redshift, Google BigQuery, etc.

Vantagens:

- Permite o armazenamento de grandes volumes de dados históricos.
- Facilita a execução de consultas analíticas complexas.
- Oferece uma base sólida para análises futuras e integração com ferramentas de BI.

## 3. Análise e Geração de Resultados com SQL

Estratégia:

- Uma vez que os dados estão no data lake, podemos utilizar a linguagem SQL para gerar queries analíticas e trazer as informações específicas que forem solicitadas.

- Falando especificamente deste case, podemos escrever uma query para obter o identificador do usuário, a data mais recente em que o usuário adicionou produtos ao carrinho e a categoria em que o usuário tem mais produtos adicionados ao carrinho.
- Esta query pode ser usada em ferramentas de BI como Tableau, Power BI, Looker, etc., permitindo ao time de data analytics visualizar e explorar os dados de maneira mais eficiente além do processo como um todo ser mais eficiente.

## Implementação Detalhada

### Código Melhorado com Paralelização e Armazenamento Incremental

[#### script melhorado ####](#)

#### Query SQL Analítica

```
SELECT
  user_id,
  MAX(last_added_date) AS last_added_date,
  top_category
FROM user_cart_data
GROUP BY user_id, top_category
ORDER BY last_added_date DESC;
```

## Estratégias Avançadas

1. **Automatização do Processo ETL:**
  - Agendar o script ETL para rodar periodicamente (por exemplo, utilizando ferramentas de orquestração de dados, como o Apache Airflow).
  - Usar serviços como AWS Lambda, Google Cloud Functions para automatizar a execução do script em resposta a eventos.
2. **Monitoramento e Logging:**
  - Implementar logging adequado para monitorar a execução do ETL e capturar quaisquer erros ou anomalias.
  - Utilizar ferramentas de monitoramento como AWS CloudWatch, Google Cloud Monitoring para alertas em tempo real.
3. **Validação de Dados:**
  - Adicionar validação de dados para garantir a integridade e qualidade dos dados carregados no data lake.
  - Implementar testes unitários para as funções de extração, transformação e carga.
4. **Integração com Ferramentas de BI:**
  - Conectar o banco de dados (data lake) com ferramentas de BI como Tableau, Power BI ou Looker para visualização e análise dos dados.
  - Tendo essa integração, o time de data analytics pode desenvolver dashboards interativos onde será possível explorar os dados facilmente.

## **Conclusão**

Estas melhorias não apenas otimizam o processo de ETL atual, mas também estabelecem uma base sólida para análises avançadas e tomadas de decisão informadas. A combinação de Python para o ETL, um data lake para armazenamento de dados históricos, e SQL para consultas analíticas proporciona uma solução robusta e escalável para gerenciamento correto e eficiente de dados.