bigDGeneral_IPCA Project

DESCRIÇÃO:

O projeto bigDGeneral_IPCA tem como objetivo baixar todo mês a tabela de IPCA do endereço https://www.idinheiro.com.br/tabelas/tabela-ipca/, extrair os dados mensais de IPCA, transformar em tabela estruturada e gravar em um arquivo parquet, fazer a ingestão em um bucket s3 na AWS para posterior catalogação desse arquivo, através do AWS Glue Data Catalog, sendo os dados posteriormente liberados para os Data Analysts e Data Scientysts que usarão ferramentas de insights para geração de apresentações gráficas para a empresa.

OBJETIVO:

Motivação:

Teve como motivador para que esse projeto bigDGeneral_IPCA fosse feito: A necessidade de se ter o índice de inflação no Data Lake bigDVarejo e com isso conhecer a variação de preços praticados nas vendas ao consumidor.

Aplicação Prática:

Conhecer qual foi a inflação de determinado período relacionada aos preços de produtos voltados ao consumidor final, ao consumidor do varejo. Através desse conhecimento usar esses índices para se descontar o que foi perdido para a inflação, fazendo o alinhamento de qualquer valor no cálculo de sua variação.

Resultados Esperados.

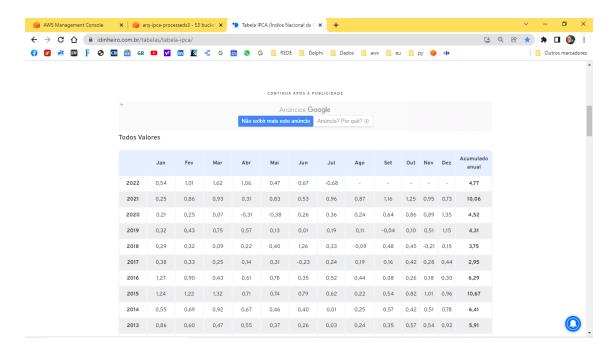
Permitir a correção de qualquer valor ou índice descontando a variação da inflação do valor de produtos voltados ao consumidor.

MÉTODOS:

De acordo com a motivação já descrita anteriormente deu-se o início para que os primeiros passos fossem dados, conforme o problema apresentado, em direção ao estudo e implementação do bigDGeneral_IPCA:

- Verificar se as solicitações apresentadas são viáveis.
 - De acordo com as necessidades apresentadas fez-se uma análise preliminar para se certificar da viabilidade do projeto.
- Pesquisar tabelas de dados relativos à variação do IPCA:
 - Procurou-se encontrar na web uma tabela que resumisse todos os índices de IPCA existentes e atualizados do Brasil.
- Avaliar arquivos de dados disponibilizados:

analisados e validados se poderiam entregar o resultado de informação que era necessário para a comparação da perda de inflação ocorrida no período.



- Verificar a viabilidade de baixar essa tabela via python:
 - Outra verificação importante foi a de se certificar que a tabela eleita para o projeto pudesse ser baixada, utilizando funções python que fazem web scraping e que estivessem disponíveis continuamente, inclusive com as atualizações.
- Definir quais os formatos de arquivos serão utilizados no armazenamento:
 - A ideia é armazenar a tabela de IPCA mensal em arquivo parquet, que será catalogado pelo Glue Data Catalog, podendo os dados e informações serem disponibilizados para o usuário final data analysts, data scientists ou mesmo data engineers.
- · Definir onde os arquivos serão armazenados, em local ou cloud:
 - O local onde os arquivos ficarão armazenados tem como peso principal a continuidade do armazenamento com boa performance de acesso. Elegeu-se o aws S3 para armazenamento, que oferece toda a manutenção da estrutura, com armazenamento distribuído e alta escalabilidade, deixando os engenheiros de dados livres para se preocuparem somente com o planejamento e operação do ELT.
- Definir quantos S3 bucket serão criados para o armazenamento:
 - Ficou definido que será necessário somente um bucket s3 para o Data Lake com o nome de arq-ipca-processeds3.
- Definir quantas camadas serão necessárias para o processamento.

- O processamento será feito em 2 camadas, uma que usará o python fazendo a extração com web scraping da tabela de IPCA gravando em arquivo parquet e uma segunda camada que fará a catalogação dos dados através do aws Glue com suas ferramentas.
- Definir qual a linguagem será utilizada para o processamento dos dados:
 - O Por conter uma enorme diversidade de bibliotecas para inúmeros fins, apresentar uma simplicidade de estrutura voltada para orientação a objetos, por ser uma linguagem que está sendo muito utilizada no mundo sendo uma tendência em manipulação de dados, apresentando funções voltadas para tal, por ser de fácil uso dentro da aws, optou-se em utilizar a linguagem python.
- · Definir onde serão executados os códigos.
 - Os códigos serão processados utilizando-se o serviço da aws Lambda, que é orientada a eventos com computação sem servidor, não é necessário definir um servidor para executar uma aplicação ficando transparente para nosso processamento, sendo mais uma preocupação para a equipe da aws Amazon manter o serviço funcionando com escalabilidade.
- Definir onde será o repositório de hospedagem dos códigos.
 - Devido à experiência com o Git-Hub, por apresentar seus recursos de hospedagem e manutenção de versões de código com simplicidade e objetividade, pela sua divulgação e utilização na comunidade de desenvolvimento de software, definiu-se pela utilização dessa plataforma.
- · Definir as bibliotecas utilizadas.
 - Para fazer upload dos arquivos baixados e gerados utilizar-se-á a biblioteca python boto3 (facilita o acesso aos serviços da aws).
- Definir as características e configurações do scheduler:
 - o O processamento das 2 camadas ocorrerá no quinto dia do mês às 24:00.

PRODUTOS, SERVIÇOS E SISTEMAS:



ESCOPO DA ARQUITETURA:

dgIPCAproject

