

bigDVarejoPMC Project

DESCRIÇÃO:

O projeto bigDVarejoPMC, PMC significa Pesquisa Mensal de Comércio, diz respeito a uma parte de um projeto maior o bigDVarejo, um trabalho que será composto por inúmeros projetos menores que alimentarão um repositório de Data Lake voltado para comportar dados do varejo. Esses dados poderão ser de diferentes origens, tanto gerados pela empresa quanto de ambientes externos a ela, permitindo ser tratados e armazenados de forma que ofereçam informações úteis e que possam ser usadas no apoio à tomada de decisões.

OBJETIVO:

Motivação:

O que levou a planejar o bigDVarejoPMC inicialmente foi a necessidade de revelar qual a situação da empresa com relação a evolução de vendas do varejo comparada às vendas do país e para isso descobriu-se a oferta de arquivos, disponibilizados pelo IBGE, que entregam dados dos percentuais de crescimento de vendas no varejo agrupados por UF (unidade da federação) como também por categoria de comércio e de produto. Entende-se aqui que se uma empresa tem um crescimento de vendas aquém do crescimento ocorrido no seu meio ambiente de atuação, algo está errado e deve ser identificado e revisto para a devida correção do problema. Essa necessidade de entender melhor o que está acontecendo e o que pode acontecer também motiva a encontrar o problema o que será possível através do cruzamento de outros dados que ajudarão na formação do Data Lake.

Aplicação Prática:

Os valores percentuais de crescimento do PMC revelam uma posição passada de como se comportou o crescimento de vendas por UF, por categoria de comércio e por categoria de produtos no Brasil, dessa forma, pode-se saber se a empresa esteve participando desse crescimento de vendas de varejo ou se esteve deficitária no período em questão. Se comparado a evolução dos percentuais de crescimento de vendas da empresa, nos períodos passados, com os percentuais do Brasil, traçando uma curva de desempenho tanto da empresa quanto do país, e ainda levando em consideração os números por estado e categorias comerciais e de produtos pode-se verificar o sucesso e ou fracasso desse crescimento de vendas no varejo.

Resultados Esperados:

A evolução dos valores percentuais de vendas no varejo oferecidas pelo IBGE torna-se aqui uma informação importante para balizar se a empresa está desempenhando um papel positivo ou negativo frente ao crescimento de vendas no varejo do país, podendo, dessa forma, oferecer um diferencial para que a equipe de data analytics possa defender sua tese sobre a posição da organização frente a esses importantes dados estratégicos. As informações aqui em questão estão apoiadas em dados oficiais do comércio varejista e disponibilizados pelo IBGE. São dados que visam fornecer aos inúmeros interessados como evoluiu o crescimento de vendas no varejo durante os períodos passados fornecidos. No caso do Data Lake do bigDVarejo, tratado

nesse projeto, tem-se o foco no desempenho da empresa em relação ao crescimento das vendas varejistas dentro do todo que é o Brasil. Espera-se que essa comparação do crescimento de vendas revele se a empresa teve uma participação positiva, participando desse crescimento ou negativa, ficando à margem da oportunidade de crescer em suas vendas no varejo, e ainda, com esse diagnóstico descritivo permitir que a equipe de Analytics possa buscar mais informações para enriquecer sua tese apresentada em insights que serão divulgadas para a gestão estratégica da companhia.

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 bigDVarejoPMC:

- *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 dados relativo à evolução de vendas do varejo no Brasil:*
 - Conforme o problema apresentado procurou-se, através da máquina de busca do google, por arquivos que poderiam conter dados para atender as necessidades apresentadas pela empresa.
Os dados mais viáveis foram encontrados no site do IBGE (<https://www.ibge.gov.br/estatisticas/economicas/comercio/9227-pesquisa-mensal-de-comercio.html?=&t=downloads>), correspondendo a tabelas de percentual de crescimento de vendas agrupados por UF, categoria de comércio e de produto.
- *Avaliar arquivos de dados disponibilizados pelo IBGE:*
 - Os dados oferecidos pelo site do IBGE foram analisados e validados se poderiam entregar o resultado de informação que era necessário para a comparação das vendas da companhia com as vendas do país.
- *Verificar a viabilidade de baixar esses arquivos via python:*
 - Outra verificação importante foi a de se certificar que os arquivos eleitos para o projeto pudessem ser baixados através de API e que estivessem disponíveis continuamente, inclusive com as atualizações.
- *Definir quais os formatos de arquivos serão utilizados no armazenamento:*
 - Os arquivos originais com os dados brutos/crus (raw) são arquivos .xls, planilhas excel baixadas do site do IBGE, que serão transformadas em arquivos .csv e posteriormente em arquivos parquet. A ideia é armazenar os arquivos brutos, em .csv, para uma possível revisão de estrutura e os arquivos parquet com os 2 tipos de tabelas separadas, UF e Categorias de comércio e produto, que serão catalogados 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. Se tivesse a estrutura montada na empresa, com profissionais disponíveis para a manutenção dessa estrutura, poder-se-ia fazer uso da mesma, mas como está-se começando a montar o Data Lake e a empresa não tem nada disso ainda, optou-se por terceirizar a estrutura de armazenamento como também algumas ferramentas de acesso aos dados, assim elegeram-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:*
 - Definido que o armazenamento dos arquivos ficará no S3 bucket e como cada bucket é praticamente uma pasta, um diretório, buscando organizar os arquivos por tipo, ficou definido que o Data Lake terá, inicialmente 3 buckets sendo definidos os nomes: arquivosPMCCrawS3 (dados brutos .xls), arquivosPMCprocessedS3 (dados transformados json) e arquivosPMCCuratedS3 (tabelas hive/parquet).
- *Definir quantas camadas serão necessárias para o processamento:*

O processamento será feito em 3 partes distintas ou 3 camadas sequenciais que são:

 - *pmcBRONZE:*

Baixa os arquivos pmc.xls (pesquisa mensal de comércio) contendo o % de crescimento das vendas por categoria comercial e UF. Transforma esse arquivo para .csv e faz um load desse arquivo para o bucket arquivosPMCCrawS3, onde ficarão armazenados para uma possível reestruturação dos dados e novo processamento.
 - *pmcSILVER:*

Transforma os arquivos .csv que estão no bucket arquivosPMCCrawS3, separa as informações de UF das informações de categoria comercial e grava em arquivos parquet em 2 tipos de tabelas uma de UF e outra de categoria comercial e de produto. Faz o load dos arquivos parquet para o bucket arquivosPMCprocessedS3 para posterior catalogação.
 - *pmcGOLD:*

Os dados gravados em parquet no bucket arquivosPMCprocessedS3 agora são catalogados pelo Glue Data Catalog ficando acessíveis pelas ferramentas de pesquisas, no caso o aws athena e aws quickinsight.

- *Definir qual a linguagem será utilizada para o processamento dos dados:*

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:*

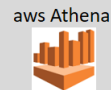
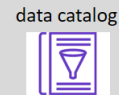
- Para fazer upload dos arquivos baixados e gerados utilizar-se-á a biblioteca python boto3 (facilita o acesso aos serviços da aws).

- *Definir as características e configurações do scheduler:*

- O processamento das 3 camadas serão executados às quartas e sextas-feiras às 20:00.

PRODUTOS, SERVIÇOS E SISTEMAS:

bigDVarejoPMC



ESCOPO DA ARQUITETURA:

