

## PROVA PRÁTICA ESPECIALISTA EM ENGENHARIA DE DADOS - DEVOPS

### OBSERVATÓRIO DA INDÚSTRIA

#### O que queremos receber?

Um repositório no seu Github com texto da pergunta 1 e 2ª, além dos scripts das demais questões. Além do repositório, queremos também receber um arquivo zip no e-mails: [dgasilva@sfiec.org.br](mailto:dgasilva@sfiec.org.br), [elgomes@sfiec.org.br](mailto:elgomes@sfiec.org.br) e [rpadilha@sfiec.org.br](mailto:rpadilha@sfiec.org.br).

#### Como a prova de ser feita?

Utilize todos os seus conhecimentos para propor a melhor solução possível para os problemas apresentados. Não se acanhe, use sua criatividade e demonstre todo o seu poder de desenvolvedor. Nosso time irá avaliar a sua capacidade através do código que será entregue, por tanto, capriche, pois seus concorrentes irão caprichar. Ser ousado não tira pontos, pelo contrário, ajuda nossa equipe a avaliar suas habilidades.

#### 1) Auto avaliação

Autoavalie suas habilidades nos requisitos de acordo com os níveis especificados.

Qual o seu nível de domínio nas técnicas/ferramentas listadas abaixo, onde:

- 0, 1, 2 - não tem conhecimento e experiência;
- 3, 4, 5 - conhece a técnica e tem pouca experiência;
- 6 - domina a técnica e já desenvolveu vários projetos utilizando-a.

#### **Tópicos de Conhecimento:**

- Manipulação e tratamento de dados com Python: \_\_\_\_
- Manipulação e tratamento de dados com Pyspark: \_\_\_\_
- Desenvolvimento de data workflows em Ambiente Azure com databricks: \_\_\_\_
- Desenvolvimento de data workflows com Airflow: \_\_\_\_
- Manipulação de bases de dados NoSQL: \_\_\_\_
- Web crawling e web scraping para mineração de dados: \_\_\_\_
- Construção de APIs: REST, SOAP e Microservices: \_\_\_\_

#### 2) **Desenvolvimento de pipelines de ETL de dados com Python, Apache Airflow, Hadoop e Spark**

Foi solicitado à equipe de AI+Analytics do Observatório da Indústria/FIEC, um projeto envolvendo os dados do Anuário Estatísticos da ANTAQ (Agência Nacional de Transportes Aquáticos).

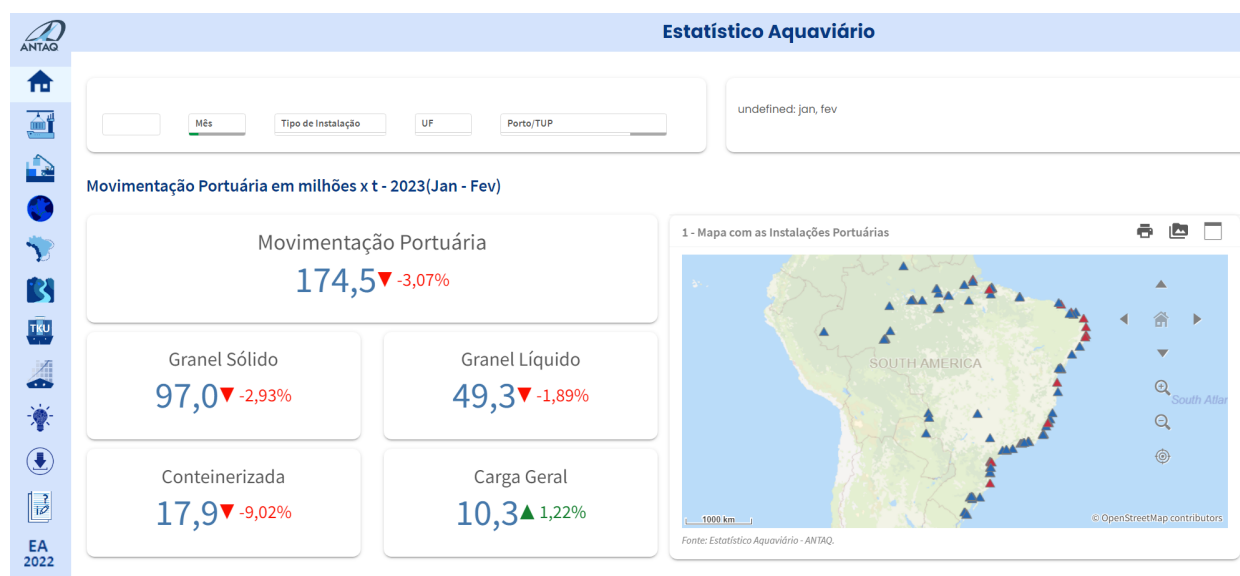
O projeto consiste em uma análise pela equipe de cientistas de dados, bem como a disponibilização dos dados para o cliente que possui uma equipe de analistas própria que utiliza a ferramenta de BI (*business intelligence*) da Microsoft.

Para isto, o nosso *cientista de dados* tem que entender a forma de apresentação dos

dados pela ANTAQ e assim, fazer o ETL dos dados e os disponibilizar no nosso data lake para ser consumido pelo time de cientistas de dados, como também, elaborar uma forma de entregar os dados tratados ao time de analistas do cliente da melhor forma possível.

### Informações Importantes:

Painel de BI: <https://web3.antaq.gov.br/ea/sense/index.html#pt>



Documentação: <https://web3.antaq.gov.br/ea/sense/download.html#pt>

**Estatístico Aquaviário**

**Base de Dados**

A seguir, apresenta-se a modelagem de dados do Estatístico Aquaviário, inclusive os relacionamentos entre as tabelas da base de dados. Observações: 1) Os dicionários de dados estão associados a cada tabela; 2) Os arquivos estão no formato de texto (txt) com separador ";".

► Siga os seguintes passos para o download:

A) Primeiro escolha o ano desejado do download, a partir do campo "Escolher Ano";

B) Opção de download de apenas uma tabela: Clique em cima da tabela desejada;

C) Opção de download de todas as tabelas: Clique em cima do ícone "Download de todos os arquivos do ano escolhido".

**Download dos Arquivos Compactados**

Escolher Ano: [2023 ▼]

**Tabela com os arquivos compactados (zip) por ano**

#	Arquivo	Download
1	Atracação	Clique aqui.
2	Carga	Clique aqui.
3	Carga Containerizada	Clique aqui.
4	Tempos Atracação	Clique aqui.
5	Taxa Ocupação	Clique aqui.
6	Carga Região Hidrográfica, Hidrovia e Rio	Clique aqui.
7	Todos os Arquivos	Clique aqui.

**Tabelas de Cadastro:**

Instalação Portuária Origem;  
Instalação Portuária Destino;  
Mercadorias;  
Mercadorias - Contêiner.

**Modelo de Dados e Dicionário de Dados**

Download modelagem dos dados.  
Download arquivo com os metadados (dicionário) dos dados.

**Banco SQL da FIEC:** SQL Server

**Banco NoSQL da FIEC:** Mongo DB

**Ferramenta dos analistas de BI do cliente:** Power BI

**Supondo que você seja nosso Especialista de dados:**

a) Olhando para todos os dados disponíveis na fonte citada acima, em qual estrutura de banco de dados você orienta guardá-los no nosso Data Lake? SQL ou NoSQL? Discorra sobre sua orientação. (1 pts)

b) Nosso cliente estipulou que necessita de informações apenas sobre as atracções e cargas contidas nessas atracções dos últimos 3 anos (2021-2023). Logo, o time de especialistas de dados, em conjunto com você, analisaram e decidiram que os dados devem constar no data lake do observatório e em duas tabelas do SQL Server, uma para atracção e outra para carga.

Assim, desenvolva script(s) em Python e Spark que extraia os dados do anuário, transforme-os e grave os dados tanto no data lake, quanto nas duas tabelas do SQL Server, sendo `atracao_fato` e `carga_fato`, com as respectivas colunas abaixo. Os scripts de criação das tabelas devem constar no código final.

Lembrando que os dados têm periodicidade mensal, então *script's* automatizados e robustos ganham pontos extras. (2 pontos + 1 ponto para solução automatizada e elegante).

Colunas da tabela `atracao_fato`:

IDAtracao	Tipo de Navegação da Atracção
CDTUP	Nacionalidade do Armador
IDBerco	FlagMCOperacaoAtracao
Berço	Terminal
Porto Atracção	Município
Apelido Instalação Portuária	UF
Complexo Portuário	SGUF
Tipo da Autoridade Portuária	Região Geográfica
Data Atracção	Nº da Capitania
Data Chegada	Nº do IMO
Data Desatracção	TEsperaAtracao
Data Início Operação	TesperalInicioOp

Data Término Operação	TOperacao
Ano da data de início da operação	TEsperaDesatracacao
Mês da data de início da operação	TAtacado
Tipo de Operação	TEstadia

Colunas da tabela carga\_fato: (atente-se que para o tipo de carga containerizada, pois cada contêiner pode ter mais de uma mercadoria)

IDCarga	FlagTransporteViaInterior
IDAtracacao	Percurso Transporte em vias Interiores
Origem	Percurso Transporte Interiores
Destino	STNaturezaCarga
CDMercadoria (Para carga containerizada informar código das mercadorias dentro do contêiner.)	STSH2
Tipo Operação da Carga	STSH4
Carga Geral Acondicionamento	Natureza da Carga
ContainerEstado	Sentido
Tipo Navegação	TEU
FlagAutorizacao	QTCarga
FlagCabotagem	VLPesoCargaBruta
FlagCabotagemMovimentacao	Ano da data de início da operação da atracação
FlagContainerTamanho	Mês da data de início da operação da atracação

FlagLongoCurso	Porto Atracação
FlagMCOperacaoCarga	SGUF
FlagOffshore	Peso líquido da carga ( Carga não containerizada = Peso bruto; Carga containerizada = Peso sem contêiner)

- c) Essas duas tabelas ficaram guardadas no nosso Banco SQL SERVER. Nossos economistas gostaram tanto dos dados novos que querem escrever uma publicação sobre eles. Mais especificamente sobre o tempo de espera dos navios para atracar. Mas eles não sabem consultar o nosso banco e apenas usam o Excel. Nesse caso, pediram a você para criar uma consulta (query) otimizada em sql em que eles vão rodar no excel e por isso precisa ter o menor número de linhas possível para não travar o programa. Eles querem uma tabela com dados do Ceará, Nordeste e Brasil contendo número de atracações, para cada localidade, bem como tempo de espera para atracar e tempo atracado por meses nos anos de 2021 e 2023. Segundo tabela abaixo: (2pts)

Localidade	Número de Atracações	Variação do número de atracação em relação ao mesmo mês do ano anterior - Bônus	Tempo de espera médio	Tempo atracado médio	Mês	Ano

### 3) Criação de ambiente de desenvolvimento com Linux e Docker.

Finalmente, este processo deverá ser automatizado usando a ferramenta de orquestração de *workflow* Apache Airflow + Docker. Escreva uma DAG para a base ANTAQ levando em conta as características e etapas de ETL para esta base de dados considerando os repositórios de data lake e banco de dados. Esta também deve conter operadores para enviar avisos por e-mail quando necessário (e.g.: caso os dados não sejam encontrados, quando o processo for finalizado, etc). Todos os passos do processo ETL devem ser listados como *tasks* e orquestrados de forma otimizada, porém não é necessário migrar o código criado anteriormente para dentro das *tasks do Airflow* (foque em mostrar o fluxo de *tasks* e as estruturas básicas de uma DAG). Caso isso seja feito, será considerado um extra. (2 pts + 1 pts)