



Avaliação de conhecimentos

Engenharia de Dados & DataOps

Nome:_____ Data:___/___/___



Instruções gerais:

A avaliação está dividida em 4 partes:

- **Manipulação e análise de dados**
- **Lógica e Programação**
- **Tratamento de dados**
- **Inglês**

A avaliação foi dimensionada para ser resolvida entre 60-120 minutos, e servirá para conhecermos seu nível de conhecimento em Engenharia de Dados e DataOps, e não é eliminatório durante o processo. Então fique tranquilo(a), use todo seu conhecimento, justifique bem as respostas, e tente resolver tudo, mesmo se não souber.... esse ponto é bem importante na avaliação. Lembre-se de ser sucinto e direto nas respostas.

Uso de tradutores ou softwares que forneçam respostas automáticas serão consideradas como plágio e, uma vez detectado, o candidato automaticamente será eliminado do processo seletivo.

A devolução deve ocorrer pelo e-mail **recrutamento@maxxidata.com**, sendo enviada até 24 horas após seu envio. Devolva a avaliação digitalizada (anexada no e-mail).

Orientações:

- As questões 1 a 3 devem ser solucionadas utilizando a linguagem de programação de sua preferência. Ou seja, não serão aceitas manipulações manuais ou análises via excel.
- Para as questões 1,2 e 3 o código pode ser disponibilizado através do link do github pessoal, descrito neste arquivo ou em um script anexado na entrega.
- Para as questões 1,2 e 3 todo o código desenvolvido deve estar comentado com as explicações de cada operação desenvolvida
- Nas questões 2 e 3 envie também o print do resultado das questões

Boa sorte!

Queremos você em nosso time Maxxidata.

Parte 1 – Manipulação e análise de dados

A partir das tabelas *cadastro.csv* e *vendas.csv*, resolva as seguintes questões:

- Qual a pessoa (cod_cadastro) que gastou mais? e a que gastou menos?
- Qual a região de procedência que gasta mais?
- Qual o produto que mais é vendido em quantidade?
- Existe alguma característica da que identifique o grupo das pessoas (top 5) que comprem mais produtos em quantidade? E as que gastam mais dinheiro?
- Você consegue ver alguma relação entre o grau de instrução e o número de filhos? Descreva a sua interpretação do caso, utilizando os dados para justificar as suas conclusões.
- Imagine que a empresa em questão queira desenvolver um programa de cashback de benefícios baseada no perfil dos clientes, sendo assim, o primeiro passo seria desenvolver este perfil. Proponha um critério para classificar os clientes em clientes *diamante*, *ouro* e *prata*. Justifique a sua resposta.

Parte 2 – Lógica e Programação – Lista reversa

Descrição

Uma lista encadeada (= linked list = lista ligada) é uma sequência de nós; cada nó contém um objeto (todos os objetos são do mesmo tipo) e o endereço do nó seguinte.

Exemplo: Laranja -> Maçã -> Banana -> Melancia

Os objetos são do tipo string e as setas representam o endereço (também chamado de ponteiro) do nó seguinte.

Lista Encadeada Reversa: Dado o ponteiro para nó de uma lista encadeada, altere os próximos ponteiros dos nós para que sua ordem seja invertida. O ponteiro fornecido pode ser nulo, o que significa que a lista inicial está vazia.

Exemplo:

llist - **1 -> 2 -> 3 -> NULL**

llist reversa - **3 -> 2 -> 1 -> NULL**

Exercícios

- Desenvolva uma função que receba o número de elementos na lista encadeada e a lista encadeada, e retorne a lista encadeada reversa. Teste a sua função para os seguintes inputs:
 - Input1: (4, [85,15,4,20])
 - Input2: (7, [1,2,3,5,6,7,8])
- Escreva um texto explicando o funcionamento do algoritmo.
- Descreva um problema do dia a dia que este algoritmo pudesse ser utilizado.

Parte 3 – Tratamento de Dados

Imagine que dentro do projeto você foi encarregado de tratar a base cadastral abaixo:

<i>nome</i>	<i>idade</i>	<i>e-mail</i>
Carolina da Silva	35	carolina.sil@gmail.com
gabrielly rezenDe	43	gabi_rezende@google.com
Bernardo Dias	70	bernardodias@aws.com
helenA gonçAlves	23	helenA.gonçAlves@gmail.com
Renan Caldeira	61	renan_calDEIRA@google.com

Lucca Lopes	46	lucca.com.br
André Ribeiro	39	andré_ribeiro.outlook.com.br
João Marcelo	32	<u>joão_ma@outlook.br</u>
Joana da Luz	22	<u>joluz@google.com.br</u>
Vicente Moreira	19	<u>vicente.moreira@maxxi.com</u>
Camila d' Rocha	51	<u>mlarocha@gmail.com</u>
Elisa Monteiro	30	<u>Elisa_mo23@yahoo.com.br</u>
Daniela Dias	43	<u>danidias@maxxi.com</u>
daniel Gomes	20	<u>danieldias@aws.com</u>
Heitor Nogueira	64	<u>nogueira_heitor@gmail.com</u>

Desenvolva um projeto que realize os seguintes tratamentos:

- i) Tratamento do campo nome
 - a. remover os caracteres especiais
 - b. remover os espaços do começo e do final da string
 - c. transformar as string em minúsculo
- ii) Tratamento do campo e-mail
 - a. Aplicar os mesmos tratamentos (a,b e c) que foram aplicados no campo nome no campo e-mail.
 - b. Criar uma coluna com valor booleano chamada *valido*, sendo 1 = e-mail válido e 0 = e-mail não válido. Os e-mails válidos são aqueles que possuem @ e terminam com “.com” ou “.com.br”

Parte 4 – Inglês

Leia o texto e responda às perguntas em inglês (a resposta deve ser redigida em inglês).

What is a data lake?

A data lake is a central storage repository that holds big data from many sources in a raw, granular format. It can store structured, semi-structured, or unstructured data, which means data can be kept in a more flexible format for future use. When storing data, a data lake associates it with identifiers and metadata tags for faster retrieval.

Coined by James Dixon, CTO of Pentaho, the term “data lake” refers to the ad hoc nature of data in a data lake, as opposed to the clean and processed data stored in traditional data warehouse systems.

Data lakes are usually configured on a cluster of inexpensive and scalable commodity hardware. This allows data to be dumped in the lake in case there is a need for it later without having to worry about storage capacity. The clusters could either exist on-premises or in the cloud.

Data lakes are easily confused with data warehouses, but feature some distinct differences that can offer big benefits to the right organizations especially as big data and big data processes continue to migrate from on-premises to the cloud. They are similar in their basic purpose and objective, which make them easily confused:

- Both are storage repositories that consolidate the various data stores in an organization.
- The objective of both is to create a one-stop data store that will feed into various applications.

However, there are fundamental distinctions between the two that make them suitable for different scenarios.

- Schema-on-read vs schema-on-write — The schema of a data warehouse is defined and structured before storage (schema is applied while writing data). A data lake, in contrast, has no predefined schema, which allows it to store data in its native format. So in a data warehouse most of the data preparation usually happens before processing. In a data lake, it happens later, when the data is actually being used.
- Complex vs simple user accessibility — As data is not organized in a simplified form before storage, a data lake often needs an expert with a thorough understanding of the various kinds of data and their relationships, to read through it. A data warehouse, in contrast, is easily accessible to both tech and non-tech users due its well-defined and documented schema. Even a new member on the team can begin to use a warehouse quickly.
- Flexibility vs rigidity — With a data warehouse, not only does it take time to define the schema at first, it also takes considerable resources to modify it when requirements change in the future. However, data lakes can adapt to changes easily. Also, as the need for storage capacity increases, it is easier to scale the servers on a data lake cluster.

For more on this distinction, and to help determine which is best for your organization, see “Data Lakes vs Data Warehouses”. There is also an emerging open data management architecture that combines the flexibility of a data lake with the data management capabilities of a data warehouse, known as a data lakehouse.

Questions

- i) In your opinion what would be the best data storage (data lake or data warehouse) to store data from different formats (files, images, etc...) and sources (data bases, api, sensors, etc..)? Justify your answer.
- ii) For a non-technical user what would be the best data storage (data lake or data warehouse) to adopt? Justify your answer.
- iii) In your opinion which data storage (data lake or data warehouse) is the most sensitive to schema changes?