



DataKnow

Prueba Técnica de Conocimiento

Perfil Analítico – Ingeniero de Datos – Científico de Datos

Muchas gracias por tu interés en participar en la convocatoria para pertenecer a la familia DataKnow. **AI Team DataKowner**

Estamos buscando personas comprometidas, que se destaquen por realizar un trabajo de calidad, con buena actitud de servicio, compromiso y mucha responsabilidad con sus actividades día a día, siempre dando prioridad a las necesidades del cliente. Además, también estamos buscando personas con ganas de adquirir sólidos conocimientos en técnicas de modelación, con habilidades en estadística, matemáticas, bases de datos, big data, nube (Azure, AWS, GCP), con fuertes habilidades en programación en lenguajes como: R, Python, SAS, PL/SQL, Scala, Hadoop, entre otros.

El propósito de esta prueba es medir sus capacidades para manipular datos de diferentes industrias, realizar supuestos, filtrar y utilizar información relevante, concluir y comunicar adecuadamente los resultados de los modelos. Pruebe usar cualquier herramienta de programación.

1. CARGA DE INFORMACIÓN

Cargar un data set, realizar el cargue y depuración del archivo OFEI1204.txt.

Se debe entregar una tabla con las columnas:

Agente

Planta

Hora_1

Hora_2

Hora_3

...

Hora_24

Solamente procesar los registros Tipo D.

Enviar junto con la tabla resultante el código utilizado.

Explicar el paso a paso en un archivo de texto (.doc o .pdf).

2. MANIPULACIÓN DE DATOS

- Cargar un data set, del archivo Excel Master Data, únicamente las siguientes columnas:
- Nombre visible Agente
- AGENTE (OFEI)
- CENTRAL (dDEC, dSEGDES, dPRU...)
- Tipo de central (Hidro, Termo, Filo, Menor)
- Seleccionar los registros que pertenecen al agente EMGESA ó EMGESA S.A. y adicionalmente que el Tipo de Central sea 'H' o 'T'.
- Cargar el archivo dDEC1204.TXT que viene por Central.
- Realizar el merge de los dos data sets por Central.
- Calcular la suma horizontal de todas las horas para cada planta.
- Seleccionar solamente los registros de las plantas cuya suma horizontal sea mayor que cero.
- Los resultados deben ser entregados en un dataset.
- Enviar junto con la tabla resultante el código utilizado.
- Explicar el paso a paso en un archivo de texto (.doc o .pdf).

3. PRUEBA DE SQL

El SQL (Structured Query Language) es un lenguaje estándar para almacenar, manipular y recuperar datos en bases de datos. Es uno de los idiomas más comunes para especificar y acceder a los datos. Responda las siguientes preguntas utilizando solo consultas SQL.

Se puede resolver usando un SQL Online. a. Explicar el paso a paso en un archivo de texto (.doc o .pdf).

Se puede resolver usando cualquier motor de base de datos o en su defecto un compilador de SQL Online como los siguientes:

<https://sqliteonline.com/>

<http://www.sqlfiddle.com/>

Nota: el código a continuación es para la creación de las tablas insumos.

```
CREATE TABLE EMPLEADO (
  ID INT(8),
  NOMBRE VARCHAR(50),
  APELLIDO VARCHAR(59),
```

```
SEXO CHAR(1),
FECHA_NACIMIENTO DATE,
SALARIO DOUBLE(10,2)
);

CREATE TABLE VACACIONES(
    ID INT(8),
    ID_EMP INT(8),
    FECHA_INICIO DATE,
    FECHA_FIN DATE,
    ESTADO CHAR(1),
    CANTIDAD_DIAS INT(8)
);

/*EN ESTA TABLA SE ALMACENA LA INFORMACIÓN BASICA DE LOS EMPLEADOS*/

INSERT INTO EMPLEADO VALUES (1,"JUAN","PELAEZ","M",'1985-01-29',3500000);
INSERT INTO EMPLEADO VALUES (2,"ANDRES","GARCIA","M",'1975-05-22',5500000);
INSERT INTO EMPLEADO VALUES (3,"LAURA","PEREZ","F",'1991-09-10',2500000);
INSERT INTO EMPLEADO VALUES (4,"PEPE","MARTINEZ","M",'1987-12-01',3800000);
INSERT INTO EMPLEADO VALUES (5,"MARGARITA","CORRALES","F",'1990-07-02',4500000);

/*EN ESTA TABLA SE ALMACENA LAS SOLCITUDES DE VACIONES DE CADA EMPLEADO*/

INSERT INTO VACACIONES VALUES (1,1,'2019-07-01','2019-07-15','A',14);
INSERT INTO VACACIONES VALUES (2,2,'2019-03-01','2019-03-15','R',14);
INSERT INTO VACACIONES VALUES (3,2,'2019-04-01','2019-04-15','A',14);
INSERT INTO VACACIONES VALUES (4,2,'2019-08-14','2019-08-20','A',6);
INSERT INTO VACACIONES VALUES (5,3,'2019-08-20','2019-08-25','A',5);
INSERT INTO VACACIONES VALUES (6,3,'2019-12-20','2019-12-31','A',11);
```

Preguntas:

- Seleccione nombre, apellido y salario de todos los empleados.
- Seleccione nombre, apellido y salario de todos los empleados que ganen más de 4 millones.
- Cuento los empleados por sexo.
- Seleccione los empleados que no han hecho solicitud de vacaciones.
- Seleccione los empleados que tengan más de una solicitud de vacaciones y muestre cuantas solicitudes tienen los que cumplen.
- Determine el salario promedio de los empleados.
- Determine la cantidad de días promedio solicitados de vacaciones por cada empleado.
- Seleccione el empleado que mayor cantidad de días de vacaciones ha solicitado, muestre el nombre, apellido y cantidad de días totales solicitados.
- Consulte la cantidad de días aprobados y rechazados por cada empleado, en caso de no tener solicitudes mostrar 0.

4. PRUEBA DE AWS

Construya una solución completa en la nube de AWS que usando todas las tablas de la fuente de datos de muestra de

Redshift: <https://docs.aws.amazon.com/redshift/latest/gsg/samples/ticketdb.zip>

Puede encontrar más información sobre los tipos y columnas de estos datos en:

Step 6: Load sample data from Amazon S3 - Amazon Redshift

Realice las siguientes tareas y guarde el código y evidencia de ejecución de las mismas:

1. Configure los roles específicos de los recursos AWS necesarios para todo el ejercicio
2. Configure un Clúster de pruebas de Redshift con los requerimientos mínimos necesarios
3. Copie la información a Redshift usando algún editor de queries como SQL Workbench/J - Home (sql-workbench.eu)
4. Responda las siguientes preguntas usando comandos de consultas SQL:
 - a. ¿Cuántos Usuarios gustan del Jazz?
 - b. ¿Cuántos Usuarios gustan de la ópera y del rock al mismo tiempo?
 - c. ¿Cuál es el promedio, moda y mediana del total de Ventas?
 - d. ¿Cuál el promedio de ventas de usuarios que gustan del rock, pero NO del Jazz?

5. En una nueva tabla junte la información (Nombre de usuario, Apellido de usuario, Correo de usuario, Nombre del evento, lugar del evento, Fecha del evento, Cantidad y Total vendidos) y expórtela usando Redshift a un bucket predefinido de S3.
6. (Opcional) Sobre la data exportada en el punto 5, y usando cualquier información adicional que desee, cree una sesión de SageMaker y realice un modelo de Forecast con cualquier técnica de su preferencia, para pronosticar las ventas para los siguientes 7 días desde el final del histórico de datos. Tenga en cuenta que la fecha de la venta se encuentra en la variable saletime y que esta está mostrada a una granularidad de factura individual.

Opcional: se puede resolver sobre tecnología Azure.

5. PRUEBA DE AZURE

Construya una solución completa en la nube de Azure que usando todas la base, de pruebas adventure Works, permita crear una etl, para la realización de un trabajo de reporteria dentro de la organización.

- desplegar base de datos en sql, con la base de pruebas adventure works
- realizar un pipeline con Azure Datafactory, utilizando data flow, para realizar la carga de una base de datos, crear 5 indicadores.
- realizar una etl, que poble un datalake.

6. PRUEBA DE MODELACIÓN ANALÍTICA

El archivo train.csv contiene información sobre muchas transacciones con tarjetas de crédito y débito por diferentes canales. Para cada transacción se tiene el valor monetario de la misma y otras variables (ver diccionario_variables.xlsx). De particular importancia es la variable FRAUDE en donde aparece 1 si la transacción constituyó un fraude o 0 si fue una transacción legítima. Su misión es desarrollar un modelo que permita, a partir de los datos en este archivo predecir cuál será el valor de la variable FRAUDE para una transacción cualquiera. El archivo test.csv contiene exactamente las mismas columnas de train.csv, la columna FRAUDE la dejamos en blanco.

1. Cargue el archivo train.csv y Construya un modelo que capaz de realizar predicciones de FRAUDE.
2. Enviar un archivo test_evaluado.csv con todas las columnas en el mismo orden que se encuentran en test.csv y adicionalmente la columna FRAUDE poblada con el valor predicho por su modelo. Cualquier valor

real (es decir, fraccionario) entre 0 y 1 será admisible aquí, donde 1 debe corresponder a FRAUDE y 0 a transacción legítima.

3. *opcional, realizar montaje de este sobre servicios azure o AWS, realizar puesta en productivo batch o/y servicio.

Nota: Muy importante enviar un archivo de texto (.doc o .pdf) donde se documente muy bien cada paso realizado, se muestren claramente los resultados y análisis respectivos, adicionalmente se deben enviar todos los códigos y comandos (con comentarios) utilizados para desarrollar esta prueba.

7. **Arquitectura**

En la empresa gaseosas SA están trabajando en una solución analítica que sea capaz de procesar miles de datos de las ventas donde se describen comportamiento de compra y análisis previos hechos por vendedores a clientes con gran volumen de compra, de forma rápida y confiable mediante el uso de tecnologías Big Data de analítica, para entrenar un modelo que sea capaz de identificar los patrones de estas ventas y compararlos en tiempo real con los patrones de datos capturados de manera streaming por dispositivos implantados puntos de venta, para controlar tempranamente y evitar el desabastecimiento.

Tu tarea es realizar un correcto diseño de la arquitectura para la solución analítica que podría soportar estos requerimientos. (Ilustra tu diseño y da una breve explicación de su funcionamiento), es importa definir el gobierno de datos y modelos de acuerdo a los perfiles

MUCHAS GRACIAS, MUCHOS EXITOS!!!!!!!