

Paralelización con Dask: ¿Cómo se aplicó la programación paralela usando Dask para realizar la carga del conjunto de datos, la eliminación de registros duplicados, y los cálculos de media y desviación estándar? Explica el procedimiento y cualquier optimización que hayas implementado para mejorar la eficiencia.

- El archivo CSV se carga en un objeto `dask.dataframe.DataFrame` mediante la función `Dask dask.dataframe.read_csv()`, en lugar de cargar todo el conjunto de datos en memoria. Dask divide automáticamente el archivo en varios fragmentos (particiones), y cada fragmento se carga y procesa de forma independiente. Esto mejora la velocidad de carga de los datos, especialmente cuando se manejan grandes conjuntos de datos.
- El método `drop_duplicates()` del `Dask DataFrame` se utiliza para eliminar los registros duplicados. En cada partición, este método también se realiza en paralelo. Dask divide los datos en trozos más pequeños y realiza el procesamiento en paralelo en cada trozo. El `DataFrame` final sin registros duplicados se produce combinando los resultados.
- Las funciones `mean()` y `std()` del `DataFrame` de Dask se utilizan para calcular la media y la desviación estándar, respectivamente. Dask divide los datos en fragmentos más pequeños y ejecuta los cálculos simultáneamente en cada fragmento. Los valores de la media y la desviación estándar de todo el conjunto de datos se calculan al final combinando los resultados parciales.

Distribución de datos y trabajo: Si tuvieras que escalar este proyecto a un conjunto de datos más grande que no cabría en la memoria de una sola máquina, ¿cómo distribuirías los datos y el trabajo entre diferentes máquinas usando Dask? Describe la estrategia que usarías y por qué crees que sería efectiva.

Crear un clúster informático: Para procesar y analizar los datos, crearía un clúster informático que pueda cooperar. Podría tratarse de un clúster local de varios nodos o de un clúster en la nube que se ejecute en Google Cloud Platform o Amazon EC2.

Dividir los datos en porciones: Después de cargar todo el conjunto de datos y dividirlo automáticamente en porciones más pequeñas utilizando la función `dask.dataframe.read_csv()` de Dask. Estas particiones serán repartidas entre el cluster de máquinas disponibles por Dask.

Las operaciones DataFrame de Dask se utilizarían para realizar cálculos y transformaciones en paralelo en cada partición distribuida en el clúster. Dask se encargará de la coordinación y ejecución eficiente de las tareas en las máquinas disponibles.

Correlaciones y análisis de algoritmos en paralelo: Cuando calculaste la matriz de correlación y aplicaste el algoritmo de tu elección, ¿cómo se benefició tu análisis de la programación paralela y distribuida? Explica cómo la paralelización y la distribución del trabajo mejoraron el rendimiento de estos cálculos y cualquier desafío que hayas encontrado en el camino.

Distribución del trabajo: Dask permite distribuir la construcción del árbol de decisión entre varios ordenadores o nodos de clúster. Esto permite que cada máquina trabaje de forma independiente en una parte del conjunto de datos, dividiendo la carga de trabajo y acelerando el proceso de entrenamiento. En comparación con la ejecución del cálculo en una sola máquina, la distribución del trabajo permite utilizar el paralelismo inherente al algoritmo del árbol de decisión y reducir el tiempo de entrenamiento.

Trabajar con grandes conjuntos de datos que no caben en la memoria de una sola máquina requiere el uso de la programación paralela, para la que Dask es especialmente bueno. En lugar de tener que cargar todo el conjunto de datos en la memoria, Dask permite cargar y procesar eficazmente particiones del conjunto de datos en cada máquina. Como resultado, se ejerce menos presión sobre la memoria de la máquina y se puede utilizar un conjunto de datos mayor para entrenar árboles de decisión.

Escalabilidad: Dask permite el escalado horizontal ampliando el clúster con máquinas adicionales. Esto significa que, simplemente añadiendo más nodos al clúster, se puede mejorar la capacidad de procesamiento y el rendimiento del árbol de decisión. A medida que se añaden más recursos al sistema, la escalabilidad permite manejar conjuntos de datos cada vez mayores y realizar cálculos cada vez más complejos.