

تقرير عن الطرق الثلاثة

Chunksize;Dask;GZIP

1*Dask+ChunkSize

في هذه الطريقة يتم قراءة الملف على دفعات صغيرة مثل `chunksizes=10000` بدل تحميله دفعة واحدة.

*في الوقت أبطأ قليلاً أي تأخذ وقت مقارنة ب `Dask` لأن كل جزء يقرأ ويعالج على حدة.

*أما في الذاكرة منخفضة جداً لأن البرنامج يحتفظ فقط بجزء صغير من البيانات.

الإيجابيات:

- ✓ مثالية للملفات الكبيرة مثل من `10GB` فما فوق.
- ✓ يمكنك معالجة كل جزء تدريجياً.

السلبيات:

- ✓ لا يمكن تنفيذ بعض العمليات إلا بعد تجميع النتائج يدوياً.
- ✓ أبطأ قليلاً من القراءة الكاملة.

2*Dask DataFrame

في مكتبة `Dask` يتم قراءة الملفات الكبيرة بشكل متوازي عبر معالج `Parallel Processing` تتعامل مع الملفات كما تفعل `Pandas` ولكن بطريقة مجزأة و موزعة تلقائياً.

*في الوقت تكون سريعة حيث أنها تستخدم `Dask` المعالجة المتعددة لتسريع القراءة والمعالجة.

*في الذاكرة يكون استهلاك متوسط.

الإيجابيات:

- ✓ أسرع أداء في الملفات الكبيرة.
- ✓ يستفيد من المعالجات المتعددة.

السلبيات:

- ✓ يحتاج مكتبة dask إضافية.

3*Pandas+GZIP

يتم ضغط الملف الى صيغة CSV.gz لتقليل حجمه ثم قراءته.
*في الوقت يأخذ وقت كبير في قراءة الملف.
*في استهلاك الذاكرة يكون استهلاك من متواضع الى كبير نوعا ما.

الإيجابيات :

- ✓ يقلل حجم الملف بنسبة كبيرة.
- ✓ مناسبة لنقل والتخزين.
- ✓ أداء مستقر رغم الضغط.

سلبيات:

- ✓ وقت القراءة أطول قليلا.
- ✓ يحتاج وقت أولي لضغط الملف.