

التقرير 1

الطرق الثلاث

Compression(gzip)

يتم ضغط الملف لتقليل حجم تخزينه
ثم قراءته بعد فك الضغط.

Dask

مكتبة تدير المعالجة المتوازية
للملفات الضخمة عبر تقسيمها تلقائياً
بين الأنوية و المعالجات.

Pandas+chunksize

قراءة الملفات الكبيرة على دفعات
صغيرة.

Chunksize لتفادي استهلاك
الذاكرة دفعة واحدة

إيجابيات كل طريقة

Compression:gzip

- ❖ يقلل الحجم التخزيني بشكل كبير جدا.
- ❖ مفيد لحفظ البيانات او مشاركتها اونلاين.

Dask

- ❖ اسرع من Pandas في ملفات الضخمة حيث يستخدم تعدد الانوية.
- ❖ يمكنه التعامل مع ملفات اكبر من حجم الذاكرة.
- ❖ واجهته سهلة الاستخدام .

Pandas+chunksize

- سهولة التنفيذ.
- لا تتطلب مكتبات إضافية.
- توفر تحكما جيدا في حجم الدفعة لتقليل استهلاك الذاكرة.

سلابيات كل طريقة

Compression gzip

- ❖ ابطئ طريقة لأنها تحتاج وقت لفك الضغط.
- ❖ استهلاك عالي لذاكرة اثناء الضغط.
- ❖ غير عملي للمعالجة الفورية اثناء التحليل.

Dask

- ❖ يستهلك ذاكرة اكثر اثناء العمليات.
- ❖ يحتاج تنصيب مكتبة إضافية.
- ❖ لا يعطي نتائج فورية أي يحتاج .compute

Pandas+chunksize

- لا تستفيد من تعدد الانوية single threads
- بطئ مقارنة ب Dask عند الملفات الضخمة.

حجم الملف ب GB	الذاكرة بMB	الوقت	الطريقة
9.61	287.64	163.53 s=2.73min	Pandas+chunksize
9.61	963.24	171.49s=2.68min	Dask
2.41	9770.67	2975.66s=49.59min	Compression gzip

- ❖ من خلال الجدول نستنتج ان الطريقة الأفضل هي Pandas+chunksize
- ❖ لأنها الأقل استهلاكاً للذاكرة ونتائجها ممتازة زمنياً ومناسبة في حالات التحليل السريع أما Dask فهي مناسبة للملفات العملاقة جداً أكثر من 20 جيجا بايت او المعالجة المتوازية لأنها تستغل كل أنوية المعالج وتتعامل مع ملفات اكبر من RAM أما في حالة كنا نريد حفظ البيانات و أرشفتها فنستعمل compression gzip لأنها تصغر الحجم كثيراً ولكن غير مناسبة للتحليل الفوري .