

الگوریتم عامل دور افتاده محلی یا (LOF) در سال ۲۰۰۰ برای کشف ناهنجاری‌ها و داده‌های پرت موجود در نقاط داده با اندازه‌گیری انحراف محلی یک نقطه داده با توجه به همسایه‌های آن ارائه شد. این الگوریتم (LOF)، در برگیرنده برخی مفاهیم موجود در الگوریتم‌های OPTICS و DBSCAN مانند فاصله از مرکز (core distance) و فاصله دسترسی‌پذیری (reachability distance) است که برای تخمین چگالی محلی مورد استفاده قرار می‌گیرند.

الگوریتم عامل دور افتاده محلی بر مبنای مفهوم چگالی محلی بنا شده و در آن محلی بودن بر اساس k نزدیک‌ترین همسایگی تعیین می‌شود که فاصله آن‌ها برای تخمین چگالی مورد استفاده قرار می‌گیرد. با مقایسه چگالی محلی یک شی با چگالی‌های همسایه‌های آن می‌توان نواحی دارای چگالی مشابه و نقاطی که اساساً چگالی کمتری نسبت به همسایه‌های خود دارند را تعیین کرد.

این موارد به عنوان دور افتادگی (داده پرت) در نظر گرفته می‌شوند. چگالی محلی به وسیله فاصله معمولی که یک نقطه داده توسط همسایه‌های خود «دسترسی‌پذیر» است تخمین زده می‌شود.

بیان ریاضی

فرض می‌شود (P) فاصله شی P از k امین نزدیک‌ترین همسایه است. توجه به این نکته لازم است که مجموعه k نزدیک‌ترین همسایه‌ها شامل همه اشیا موجود در این فاصله است که در صورت وقوع «گره» ممکن است بیش از k شی باشند. مجموعه k نزدیک‌ترین همسایگی به صورت $N_k(P)$ نشان داده می‌شود.

مراحل محاسبه

۱. انتخاب پارامتر k : تعداد همسایه‌های مورد بررسی تعیین می‌شود.
۲. محاسبه فاصله k -امین همسایه: برای هر نقطه مانند P فاصله‌ای به نام k -distance(P) محاسبه می‌شود.
۳. فاصله دسترسی‌پذیری (Reachability Distance): پایدارسازی محاسبات، فاصله بین دو نقطه به صورت اصلاح شده تعریف می‌شود.

$$o \in N_k(P)$$

$$RD_k(P, o) = \max(d(P, o), k - \text{distance}(o))$$

۴. چگالی دسترسی‌پذیری محلی (LRD): برای هر نقطه p ، چگالی محلی با استفاده از فاصله‌های دسترسی‌پذیری به همسایگان محاسبه می‌شود.

$$LRD_k(p) = \left(\frac{\sum_o RD_k(p, o)}{|N_k(p)|} \right)^{-1}$$

۵. محاسبه LOF: نسبت میانگین چگالی همسایگان به چگالی نقطه مورد نظر محاسبه میگردد.

$$LOF_k(P) = \sum_o \frac{LRD_k(o)}{|N_k(P)|}$$

تفسیر مقدار LOF:

اگر $1 \approx LOF_k(p)$ آنگاه چگالی p مشابه همسایه هاست.

اگر $1 > LOF_k(p)$ کمتر از همسایه هاست.

هرچه مقدار LOF بزرگتر باشد احتمال پرت بودن قوی تر است.

مزایا:

ماهیت نسبی: داده های پرت تنها در مقایسه با همسایگان مشخص می شوند.

کاربرد در داده های چند بعدی: بدون نیاز به ساده سازی یا فروض خاص درباره توزیع داده ها

انعطاف پذیری: خوش های پرت ولی معتبر را به عنوان داده های پرت شناسایی نمی کند.

معایب:

حساسیت به انتخاب پارامتر k: انتخاب نامناسب می تواند موجب کاهش دقت شود.

پیچیدگی محاسباتی: به دلیل نیاز به محاسبه فاصله بین نقاط متعدد، هزینه زمانی بالایی دارد.

محدودیت در داده های کم تراکم: در صورت یکنواخت بودن کل داده ها، کارایی کاهش می یابد.