

تاریخ دریافت مقاله: ۹۸/۰۱/۱۱

تاریخ پذیرش مقاله: ۹۹/۰۱/۲۴

بیشینه‌یابی با استفاده از مقایسه‌های نادقیق: بررسی نظری و تجربی

محمد فرشی*

استادیار، آزمایشگاه الگوریتم‌های ترکیب‌یاتی و هندسی، دانشکده علوم ریاضی، دانشگاه یزد، یزد، ایران
پست الکترونیکی: mfarshi@yazd.ac.ir

پریا قلی‌پور لوائی

کارشناس ارشد، آزمایشگاه الگوریتم‌های ترکیب‌یاتی و هندسی، دانشکده علوم ریاضی، دانشگاه یزد، یزد، ایران
پست الکترونیکی: paria.gholipour@stu.yazd.ac.ir

چکیده

گواه بر کارایی این الگوریتم با ورودی تصادفی نسبت به داده‌های خوشه‌ای و مرتب‌شده است، نتایج نظری، برای رسیدن به خطای $k = \lceil \log \log n \rceil$ که n اندازه ورودی است، تعداد $\Omega(n^{1+1/(2^k-1)})$ مقایسه را لازم می‌داند، در حالی که بررسی این مقاله نشان می‌دهد در عمل، روی داده‌های تصادفی و روی چند توزیع دیگر، تعداد مقایسه حداکثر چند برابر اندازه ورودی، برای رسیدن به این خطا کافی است. بدین جهت در عمل الگوریتم کارایی بهتری نسبت به حالت نظری دارد. تاثیر پارامترهای مختلف بر نحوه تغییر رفتار الگوریتم در «تعداد مقایسه‌های انجام‌شده» و «خطای خروجی» نیز در این مقاله بررسی می‌شوند. قبلاً بررسی عملی این الگوریتم‌ها انجام نشده است و این مقاله اولین بررسی تجربی این الگوریتم‌ها است. نتیجه این بررسی نشان از کارایی خوب این الگوریتم در کاربرد است.

واژه‌های کلیدی: الگوریتم بیشینه‌یابی، مقایسه‌های

نادقیق، الگوریتم k-MaxFind، توزیع تصادفی

بیشینه‌یابی، یکی از مسائل پایه‌ای در علم رایانه است و الگوریتم‌های متعددی وجود دارند که با استفاده از مقایسه عناصر، نسبت به بیشینه‌یابی اقدام می‌کنند. در مدل‌های مرسوم، مقایسه دو عنصر به صورت دقیق انجام می‌شود. ما مدل مقایسه‌های نادقیق را در نظر می‌گیریم: به این معنی که، اگر تفاوت مقدار دو عنصر زیاد باشد، حاصل مقایسه دقیق است؛ در غیر این صورت ممکن است نتیجه مقایسه با واقعیت منطبق نباشد. در الگوریتم‌های موجود برای بیشینه‌یابی در مدل مقایسه‌های نادقیق، ممکن است خطای ناشی از مقایسه‌های نادقیق انباشته شود و موجب خطای زیاد در خروجی نهایی گردد. می‌توان با هزینه انجام تمام مقایسه‌ها، به پایین‌ترین حد ممکن خطا رسید. اما هدف، دستیابی به یک تعادل مطلوب بین تعداد مقایسه‌های انجام‌شده و میزان خطا است. با چنین شرایطی، فلدمن و همکارانش الگوریتم k-MaxFind را ارائه دادند. در این مقاله این الگوریتم پیاده‌سازی شده و تحلیل تجربی انجام‌شده،

* نویسنده مسئول

کار با داده‌ها در کامپیوتر از جنبه‌های مختلف با عدم دقت همراه است. از یک طرف، داده‌های مورد پردازش در رایانه توسط دستگاه‌های دیجیتالی مختلف و حسگرهای متنوع جمع‌آوری می‌شوند که همه دارای حافظه بسیار محدود هستند و لذا داده‌های جمع‌آوری شده را با تقریب ارائه می‌کنند. همچنین با توجه به دقت محدود رایانه، گاهی انجام مقایسه داده‌های بسیار نزدیک به هم که تفاوت آن‌ها در رقم‌های خارج از محدوده دقت ذخیره‌سازی داده‌ها در رایانه است نیز می‌تواند با عدم دقت همراه باشد. این عدم دقت، در صورت مدیریت نشدن، می‌تواند به خطاهای بزرگ در نتیجه محاسبات منجر شود و نتایجی کاملاً متفاوت با واقعیت ارائه کند. این موضوع در استفاده از دستگاه‌های مرتبط با رایانه که با دقت محدود هستند نیز می‌تواند اتفاق بیفتد. علاوه بر این، در بررسی‌های نظیر بررسی داده‌های اجتماعی که افراد در تکمیل پرسش‌نامه‌های آن دخیل هستند، دریافت نظرات کاربران همواره با عدم دقت است و همواره این عدم دقت در تحلیل داده‌های مربوطه مد نظر قرار می‌گیرد. این موضوع می‌تواند در مدل‌سازی سیستم‌های مختلف نیز اتفاق بیفتد و عدم توجه به وجود خطا در داده‌ها و عدم دقت محاسبات، ممکن است به مدلی با رفتار کاملاً متفاوت از سیستم اصلی منجر شود.

به عنوان مثال، در مرجع [۱] الگوریتم‌هایی تحمل‌پذیر نقص برای تراشه‌ها و شبکه‌ها ارائه شده است و میزان کارایی آن نیز با شبیه‌سازی بررسی شده است. در این کار فرض بر این است که داده‌های داده شده به الگوریتم و همچنین الگوریتم شبیه‌سازی کاملاً دقیق است و انجام محاسبات و مقایسه‌ها نیز در تمام روال این الگوریتم‌ها کاملاً دقیق انجام می‌شود. چنین شرطی در واقعیت اتفاق نمی‌افتد و چه بسا ورود یک خطای ناچیز در روند محاسبات و تجمیع این خطا در روند الگوریتم، به نتیجه‌ای کاملاً متفاوت با واقعیت منجر شود.

لذا در بحث طراحی الگوریتم، به جای استفاده از

مدل‌های کلاسیک، که فرض انجام محاسبات و مقایسه‌ها به صورت دقیق را دارند که منطبق بر واقعیت رایانه‌های موجود نیست، مدل را به این صورت تغییر داده‌اند که مقایسه داده‌های نزدیک به هم، با توجه به دقت محدود ذخیره‌سازی می‌تواند با خطا انجام شود و در طراحی الگوریتم در این مدل، به امکان وجود خطا در مقایسه‌ها توجه شده است تا از هم‌افزایی خطاهای کوچک و منجر شدن آن به خطای زیاد در جواب الگوریتم جلوگیری شود. طراحی الگوریتم‌ها در این مدل، عملاً سخت‌تر از مدل مرسوم طراحی الگوریتم است، اما به دلیل انطباق بیشتر آن با رایانه، الگوریتم‌های طراحی شده در این مدل مطمئناً کاربرد بیشتری در عمل خواهند داشت.

یائو و یائو [۱] می‌خواستند بیشینه عنصر را در یک شبکه مرتب از مقایسه‌گرها پیدا کنند. هر مقایسه‌گر یک دستگاه ۲-ورودی ۲-خروجی است که قادر به مرتب‌سازی دو عدد می‌باشد. اگر e مقایسه‌گر معیوب باشند؛ یعنی بدون انجام مقایسه، ورودی‌ها را مستقیماً به عنوان خروجی در نظر بگیرند، آنگاه برای بیشینه‌یابی، $(e+1)(n-1)$ مقایسه لازم و کافی است.

راویکومار و همکارانش [۲]، با استفاده از پرس‌وجوهای مقایسه‌ای که پاسخ «بله» یا «نه» را دریافت می‌کردند، بزرگ‌ترین عدد در مجموعه‌ای از n عدد صحیح متمایز را یافتند. آن‌ها نشان دادند، برای این بیشینه‌یابی، اگر حداکثر e پرس‌وجو پاسخ اشتباه دریافت کند، $(e+1)n-1$ مقایسه لازم و کافی است. با فرض محدودیت بیشتر، که خطا فقط به پاسخ‌های «نه» محدود شود و پاسخ‌های «بله» مطمئناً درست باشد، آن‌ها ثابت کردند $2n + 2e - 4$ مقایسه کافی است. آن‌ها با گسترش مسئله برای خطاهای دلخواه، نشان دادند که $O(en)$ مقایسه لازم و کافی است. در نهایت فلدمن و همکارانش [۳] مسئله بیشینه‌یابی در مدل مقایسه‌های نادقیق را، با ارائه الگوریتم قطعی k -MaxFind حل کردند.

در این مقاله، ابتدا الگوریتم فلدمن ارائه می‌شود و

سپس عملکرد این الگوریتم روی داده‌های تصادفی و با برخی توزیع‌های دیگر بررسی می‌شود. مطالعه تجربی عملکرد الگوریتم‌ها، یک رویکرد پژوهشی است که طی سال‌های اخیر به آن بهای بسیار داده می‌شود. دلیل این امر این است که اولاً بسیاری از الگوریتم‌های ارائه شده به صورت نظری، به قدری پیچیده است که امکان پیاده‌سازی و استفاده از آن‌ها حتی سال‌ها پس از ارائه الگوریتم وجود نداشته است و عملاً این الگوریتم‌ها در محدوده نظری باقی مانده‌اند و وارد بحث کاربرد نشده‌اند. ثانیاً در برخی الگوریتم‌ها یک شکاف عمیق بین نتایج تحلیل زمان اجرای الگوریتم و آنچه در عمل اتفاق می‌افتد وجود دارد. در برخی الگوریتم‌ها، تحلیلی نظری حاکی از سرعت اجرای بالای الگوریتم است در حالی که پس از پیاده‌سازی و مقایسه عملکرد الگوریتم با الگوریتم‌های دیگر حل مسئله، در عمل چنین کارکردی مشاهده نمی‌شود. دلایل مختلفی نیز برای این امر وجود دارد که به‌عنوان نمونه می‌توان به فاصله بین مدل تحلیل نظری الگوریتم‌ها و رایانه‌های واقعی به دلیل ساده‌سازی مدل و همچنین حذف ثابت‌ها در تحلیل نظری الگوریتم‌ها اشاره کرد.

در این مقاله، تلاش شده است که عملکرد عملی الگوریتم محاسبه بیشینه روی داده‌هایی تصادفی بررسی و میزان کارایی آن‌ها در عمل مشخص شود. بررسی تجربی ما نشان می‌دهد، تعداد مقایسه‌های لازم برای رسیدن به یک دقت مشخص، در عمل بسیار کمتر از تعداد بیان شده در نتایج عملی است. این الگوریتم روی داده‌هایی با اندازه بین ۱ میلیون و ۱۰ میلیون و با توزیع‌های یکنواخت، مرتب‌شده و خوشه‌ای بررسی شده است و تعداد مقایسه‌های انجام شده در الگوریتم و خطای حاصل شده در عمل برای این داده‌ها به دست آمده است. نتایج حاصل شده، حاکی از تفاوت بسیار زیاد بین نتایج عملی و نتایج نظری در خصوص این الگوریتم است. به‌عنوان مثال، برای یک میلیون داده با توزیع یکنواخت، برای رسیدن به خطای ۲، کران‌های نظری تعداد بیش از ۲ میلیارد مقایسه را لازم

می‌داند، در حالی که طبق بررسی تجربی انجام شده، الگوریتم‌ها با انجام حداکثر ۳ میلیون مقایسه، بیشینه را با خطای حداکثر ۱,۵ به دست می‌دهد. این رفتار در داده‌های با توزیع‌های دیگر نیز مشاهده می‌شود. با افزایش اندازه ورودی، این تفاوت نیز بیشتر می‌شود. این بررسی نشان می‌دهد که این الگوریتم می‌تواند در عمل بسیار بهتر از نتایج نظری عمل کند و می‌تواند در کاربردهای بسیاری که عملاً داده‌های ورودی طبق یک توزیع تصادفی توزیع شده‌اند، مورد استفاده موفق قرار گیرد.

طبق اطلاعات ما، تاکنون بررسی عملکرد عملی این الگوریتم‌ها صورت نگرفته است و این اولین کار در بررسی عملی این الگوریتم‌ها است. در ادامه بعد از بیان مفاهیم موردنیاز، این الگوریتم، نتایج نظری حاصل از تحلیل آن و پیاده‌سازی عملکرد الگوریتم مطرح می‌شوند.

۲- مفاهیم و اصطلاحات

فرض کنید مجموعه $X = \{x_1, x_2, \dots, x_n\}$ به‌عنوان ورودی داده شده است، به‌طوری‌که هر x_i در آن، دارای مقدار نامعلوم $val(x_i)$ باشد. x_i نشان‌دهنده عنصر نام است که گاهی اوقات، به‌عنوان مقدار آن عنصر نیز استفاده می‌شود. به‌عنوان مثال $x_i < x_j$ می‌تواند بیان‌گر $val(x_i) < val(x_j)$ نیز باشد. دقت کنید که $|X| = n$ است.

در این مقاله، ما یک مدل ساده از مقایسه‌های نادقیق را در نظر می‌گیریم. فرض کنید دو خواهیم دو عنصر x_i و x_j را باهم مقایسه کنیم. در مدل مقایسه‌های نادقیق، $\delta > 0$ وجود دارد که اگر $|val(x_i) - val(x_j)| > \delta$ ، آنگاه نتیجه مقایسه دقیق و درست است، در غیر این صورت نتیجه مقایسه، نادقیق است و نمی‌توان تشخیص داد که کدام عنصر بزرگ‌تر می‌باشد. با چنین تعریفی ممکن است پیدا کردن یک بیشینه دقیق و درست، غیرممکن باشد. فرض کنید مقادیر همه عناصر ورودی در حدود δ باشد، پس تمامی مقایسه‌ها به‌صورت نادقیق انجام می‌شوند. در نتیجه

۳- الگوریتم قطعی برای بیشینه‌یابی

فرض کنید X^* بیشینه دقیق اعداد ورودی باشد. دقت کنید اگر چندین عنصر با بیشینه مقدار وجود داشته باشند یکی از آن‌ها به دلخواه انتخاب می‌شود. هدف یافتن عنصر بیشینه، تنها با استفاده از مقایسه‌های جفت‌جفت است.

با توجه به این‌که هیچ الگوریتم قطعی بیشینه‌یابی با خطای کمتر از ۲ وجود ندارد (قضیه ۳-۲ [۱])، الگوریتم 2-MaxFind با انجام حداکثر $3n^{3/2}$ مقایسه یک عنصر با مقدار حداقل $x^* - 2$ را برمی‌گرداند (لم ۳-۳ [۱]). در نظر داشته باشید که $s > 1$ در الگوریتم یک پارامتر با مقدار دلخواه است که اگر مقدار آن مشخص نشود، به صورت پیش‌فرض $|\sqrt{n}|$ در نظر گرفته می‌شود. منظور از تورنمنت با گردش نوبت، انجام تمام $n(n-1)/2$ مقایسه، سپس مرتب‌سازی عناصر با توجه به تعداد بردهایشان به منظور یافتن برنده نهایی تورنمنت است.

الگوریتم 2-MaxFind یک عنصر بیشینه را با خطای ۲ می‌یابد. این الگوریتم را طوری تغییر می‌دهیم که به جای یک عنصر بیشینه، مجموعه‌ای شامل ۱- بیشینه را برگرداند. این الگوریتم که همان الگوریتم 1-Cover است، بجای $x^* - 2$ ، یک مجموعه ۱- بیشینه را برگرداند.

کلید گام بازگشتی از الگوریتم بیشینه‌یابی با خطای کلی، الگوریتم 1-Cover است. الگوریتم 1-Cover با تنظیم پارامتر $s = \lfloor 2\sqrt{n} \rfloor$ در الگوریتم 2-MaxFind به دست می‌آید. با این تفاوت که خروجی آن اجتماعی از x ‌ها، در گام ۳ در هر بار تکرار حلقه، به‌علاوه مجموعه از اندازه $\lceil \log(s) \rceil$ که ۱-بزرگ‌تر از عناصر نامزدهای باقی‌مانده در گام ۵ است.

الگوریتم 2-MaxFind

۱. ورودی: $X = \{x_1, x_2, \dots, x_n\}$
۲. خروجی: یک عنصر با مقدار حداقل $x^* - 2$
۳. همه x_i ‌ها را به‌عنوان نامزد برچسب‌گذاری کن.
۴. تا زمانی که بیش از یک عنصر نامزد وجود دارد مراحل ۳ و ۴ را انجام ده:

در بدترین حالت هرکدام از عناصر ورودی می‌تواند عنصر بیشینه باشد. با این حال ممکن است شناسایی یک بیشینه تقریبی، شدنی باشد. عنصر $x_i^* \in X$ یک بیشینه تقریبی، یا k -بیشینه برای همه $x_j \in X$ است، اگر بتوان آن را به‌گونه‌ای یافت که عبارت $|val(x_j) - val(x_i^*)| \leq k\delta$ برای هر $x_j \in X$ و برخی k ها (ترجیحاً کوچک و طبیعی) برقرار باشد. دقت کنید هدف اصلی برای یافتن چنین x_i^* ‌ای، به حداقل رساندن تعداد مقایسه‌های انجام‌شده است. منظور از خطا، کمینه مقدار k به صورتی است که خروجی الگوریتم در نامساوی فوق صدق کند.

بدون کاستن از کلیت مسئله، ما $\delta = 1$ را در نظر می‌گیریم. چراکه مسئله با $\delta > 0$ دلخواه، معادل با $\delta = 1$ و ارزش ورودی x_i/δ است. برای مقایسه دو عنصر x و y ، اگر مقایسه‌گر ادعا کند x از y بزرگ‌تر است؛ یعنی $x \geq y - 1$ باشد، آنگاه می‌گوییم x ، y را شکست داده است. به‌سادگی نیز می‌توانیم بگوییم y از x شکست خورده است. به‌طور مشابه x از y ، k -بزرگ‌تر است، بدین معنی که $x \geq y - k$ باشد. یک عنصر k -بیشینه از یک مجموعه است، اگر آن عنصر k -بزرگ‌تر از همه عناصر در مجموعه باشد. اگر مجموعه به صراحت مشخص نشده باشد، آنگاه ما به مجموعه‌ای از تمام عناصر ورودی مراجعه می‌کنیم. یک مجموعه، مجموعه k -بیشینه است، اگر یک زیرمجموعه از همه عناصر باشد و شامل حداقل یک k -بیشینه (یعنی عنصر از ارزش حداقل $x^* - k$) باشد.

لازم به ذکر است، اگر دو عنصر قبلاً باهم مقایسه شده باشند و در روند اجرای الگوریتم احتیاج به مقایسه مجدد همین دو عنصر باشد، مقایسه دوباره انجام می‌شود و از نتایج مقایسه انجام‌شده در مراحل قبل استفاده نمی‌شود. دقت کنید تمامی لگاریتم‌های این مقاله در پایه دو هستند. برای راحتی کار، گاهی اوقات خطای گرد کردن، در نظر گرفته نمی‌شود. اگر کف و سقف اعداد تاثیر جزئی روی کران‌ها داشته باشند، در محاسبات لحاظ نمی‌شوند.

جدول ۱: بررسی اجمالی کران‌های تعداد مقایسه‌ها در الگوریتم‌های قطعی با تنظیم $a = 1$ و خطای k

کران پایین	کران‌های بالا		
	$k = \log \log n$	k	$k = 2$
$\Omega(n^{1+(2^k-1)})$	$O(n)$	$O(n^{1+(3 \cdot 2^{k-2}-1)})$	$2n^{3/2}$

حداقل دارای خطای $O(1) - \log \log n$ است.

۴- پیاده‌سازی الگوریتم بیشینه‌یابی قطعی

با توجه به این که الگوریتم k -MaxFind تنها با استفاده از مقایسه کردن عناصر ورودی و بدون نیاز به دانش درباره چگونگی رفتار توابع مورد استفاده در این الگوریتم، یک k -بیشینه متناسب با ورودی را می‌یابد، استفاده از این الگوریتم جهت رسیدن به تعادل قابل قبول بین میزان خطا و تعداد مقایسه‌های انجام‌شده، کافی است. از همین رو و با توجه به این که الگوریتم 1-Cover یک زیرروال برای الگوریتم k -MaxFind است و 2-MaxFind، نیز بعد از کاهش اندازه ورودی فقط یکبار جهت پیدا کردن بیشینه نهایی مورد استفاده قرار می‌گیرد، پس نیازی به پیاده‌سازی جداگانه و تحلیل رفتار آن‌ها نیست. در این مقاله به‌طور خاص درباره پیاده‌سازی الگوریتم k -MaxFind فارغ از چگونگی رفتار توابع مورد استفاده در آن و تغییر رفتار این تابع نسبت به تغییر پارامترهای مختلف پرداخته می‌شود. بنابراین از این پس هر جا از کلمه «الگوریتم» استفاده شد مقصودمان الگوریتم k -MaxFind است.

در پیاده‌سازی انجام‌شده، عناصر ورودی از بازه $X = [-100000.000, +100000.000]$ به تصادف انتخاب می‌شوند. حالتی را در نظر بگیرید که دو عنصر یکسان باشند. مقایسه نادقیق این دو عنصر خطایی را ایجاد نمی‌کند. پس برای دستیابی به خطای پیاده‌سازی نزدیک واقعیت، نیاز به استفاده از اعداد تصادفی است که تکراری نباشند.

برای تحلیل بهتر الگوریتم، تمامی خروجی‌ها، میانگین ۵۰ بار اجرای الگوریتم روی ورودی‌های متفاوت است. فقط به منظور تحلیل آسان‌تر داده‌های خروجی در بخش

۵. زیرمجموعه دلخواه k تایی از عناصر نامزد را انتخاب کن و روی آن‌ها یک تورنمنت با گردش نوبت برگزار کن.

۶. x را با همه عناصر نامزد مقایسه کن و عناصری که از x باخته‌اند را حذف کن.

۷. روی عناصر نامزد باقی‌مانده (که حداکثر ۳ تا هستند)، یک تورنمنت با گردش نوبت برگزار کن و عنصر با بیشترین تعداد برد را برگردان.

الگوریتم K-MaxFind

ورودی: مجموعه $X = \{x_1, x_2, \dots, x_n\}$

خروجی: یک k -بیشینه

۱. اگر $n \leq 81$ بود، خروجی الگوریتم 2-MaxFind را روی عناصر ورودی برگردان.

۲. ابتدا n عنصر ورودی را در مجموعه‌های S_1, \dots, S_r به صورت مساوی افراز کن، به طوری که هر کدام از آن‌ها دارای اندازه $r = \max\{81, 4 \cdot n^{8/(3 \cdot 2^k - 4)}\}$ باشد.

۳. الگوریتم 1-Cover روی هر S_i به منظور بازیافتن یک مجموعه T_i ، که 1-بزرگتر از S_i باشد، فراخوانی کن. $\epsilon = U_{i=1}^r (T_i)$ را بازگردان. (این فرآیند بازگشتی در 2-MaxFind متوقف می‌شود).

۳-۱- نتایج نظری

کران‌های نظری تعداد مقایسه‌ها در الگوریتم قطعی بیشینه‌یابی، به‌طور خلاصه در جدول ۱ آورده شده است. نحوه محاسبه این کران‌ها نیز، در بخش ۳ مقاله [۱] به‌طور کامل شرح داده شده است.

نتیجه نظری ۱-۲-۱: یک الگوریتم بیشینه‌یابی با استفاده از $O(n)$ مقایسه و خطای $\log \log n$ وجود دارد. قضیه ۲-۱-۲: هر الگوریتم قطعی بیشینه‌یابی مانند A با خطای k نیاز به انجام $\Omega(n^{1+(2^k-1)})$ مقایسه دارد. پس با بیان قضیه ۱ می‌توان دریافت که نتیجه ۱ اساساً محکم است.

نتیجه نظری ۲-۱-۳: اگر A یک الگوریتم بیشینه‌یابی قطعی باشد که $O(n)$ مقایسه را انجام می‌دهد، آنگاه A

۳-۴ از ورودی حالت اول استفاده شده است. معمولاً اجرای الگوریتم روی یک مجموعه ورودی همواره یک جواب دارد. اما در مورد این الگوریتم به دلیل مقایسه نادقیق که مقدار آن تصادفی است، این اتفاق رخ نمی‌دهد. پس برای ۵۰ بار اجرای الگوریتم دو روش امکان‌پذیر است:

روش اول: آرایه‌ای از اندازه مورد نیاز با عناصر تصادفی ایجاد می‌شود. این آرایه را به منظور اجراهای مجدد الگوریتم نگهداری کنید. پس الگوریتم ۵۰ بار روی همین آرایه اجرا می‌شود.

روش دوم: مشابه روش اول آرایه‌ای از اندازه مورد نیاز با عناصر تصادفی ایجاد می‌شود. الگوریتم با چنین ورودی اجرا می‌شود. برای اجرای مجدد الگوریتم آرایه جدید با عناصر تصادفی تولید کنید. پس برای ۵۰ بار اجرا الگوریتم از ۵۰ آرایه تصادفی متفاوت استفاده می‌شود.

در اینجا برای اجرای الگوریتم، هر دو حالت ورودی مورد بررسی قرار گرفته است. همان‌طور که انتظار می‌رفت، حالت دوم به دلیل تنوع بیشتر الگوهای ورودی، موجب خطای جزئی بیشتری نسبت به حالت اول در برخی از اجراها شد. اما با توجه به انتخاب دلخواه عناصر در تورنمنت با گردش نوبت و ماهیت تصادفی بودن نتیجه در مقایسه نادقیق، به‌طور کلی تفاوت قابل ملاحظه‌ای بین دو حالت ورودی دیده نشد.

پارامترهایی که تغییر رفتار آن‌ها در این پیاده‌سازی بررسی شده‌اند: اندازه آرایه ورودی (n)، آستانه تفاوت (δ) و خطای نظری (k) هستند. s نیز تعداد عناصری است که در مراحل اجرای الگوریتم، روی آن‌ها تورنمنت با گردش نوبت برگزار می‌شود. پس s پارامتری است که تغییرات آن نیز در روند اجرای الگوریتم موثر است. نکته جالب توجه در فرایند پیاده‌سازی، تأثیری است که این پارامترها روی تعداد مقایسه‌های انجام‌شده و خطای خروجی الگوریتم دارند.

۴-۱-۱ جزئیات اجرا

الگوریتم در نرم‌افزار DEV-C++ نسخه 5.6.3 پیاده‌سازی و اجرا شده است. سپس نتایج حاصل، با استفاده از نرم‌افزار متلب نسخه R2013a در قالب نمودار ترسیم شده‌اند. جزئیات سخت‌افزاری رایانه‌ای که به وسیله آن پیاده‌سازی‌ها انجام شده است، از قرار زیر می‌باشد:

2.50GHz Processor: Intel(R) Core(TM) i5-4200M
CPU @ 2.500 GHz

Installed memory (RAM): 6 GB

System type: 64-bit Operating System, x64-based processor

Windows edition: Windows 8 Pro

۴-۲- نتایج تجربی

در تحلیل رفتار الگوریتم، دو پارامتر سنجیده خواهند شد.

اولین پارامتر، تعداد مقایسه‌های انجام‌شده در الگوریتم برای پیدا کردن یک k -بیشینه است. در شمارش تعداد مقایسه‌های انجام‌شده توسط الگوریتم، مجموع تعداد مقایسه‌های دقیق و نادقیق فارغ از معنا و مفهوم آن‌ها، در نظر گرفته خواهد شد. منظور از پیچیدگی الگوریتم نیز همین تعداد مقایسه‌های انجام‌شده توسط الگوریتم است (نه زمان اجرای الگوریتم).

دومین پارامتر خطای حاصل از پیاده‌سازی در خروجی نهایی الگوریتم، با نماد k' می‌باشد. اگر مقدار بیشینه واقعی و x^* خروجی الگوریتم باشد آنگاه مقدار k' از طریق فرمول $(x_{real} - x^*) / \delta$ محاسبه می‌شود.

در این بخش درباره جزئیات نتایج پیاده‌سازی با در نظر گرفتن این‌که الگوریتم همواره تضمین می‌کند خطای پیدا کردن یک k -بیشینه (k')، با انواع توزیع‌های داده‌ای، در عمل به مراتب از خطای به‌دست آمده از نتایج نظری (یعنی k) کمتر است، می‌پردازیم.

در پیاده‌سازی، ورودی الگوریتم داده‌های علامت‌دار با ۳ رقم اعشاری است که از بازه $X = [-100000.000, +100000.000]$ انتخاب می‌شود تا بتوان با ایجاد الگوهای متنوع و محدود به بازه مشخص، جزئیات تغییر رفتار الگوریتم نسبت به

جدول ۲: نتایج حاصل از اجرای الگوریتم با ورودی با توزیع تصادفی، مرتب‌شده و خوشه‌ای با تنظیم $\delta = 1$ و $k = \lceil \log \log n \rceil$

تعداد مقایسه‌های انجام شده			خطا یا k'			n	
max	Aver	min	max	aver	min		
۲۸۶۹۱۹۲	۲۸۶۴۴۸۸	۲۸۵۹۶۵۶	۱,۴۹	۰,۴۵۹۸	۰	1×10^6	توزیع تصادفی
۱۷۱۹۹۹۷۳	۱۷۱۸۸۰۸۰	۱۷۱۷۶۲۴۴	۱,۲۹	۰,۵۳۱۲	۰	6×10^6	
۲۸۶۶۳۹۹۰	۲۸۶۵۱۱۶۶	۲۸۶۳۳۲۱۷	۱,۱۹۳	۰,۴۲۷۸	۰	1×10^7	
۲۹۷۵۵۳۰	۲۹۶۸۰۴۶	۲۹۵۹۴۹۴	۱,۳۲۱	۰,۴۱۸۵	۰	1×10^6	توزیع مرتب‌شده
۲۶۹۹۳۳۰۰	۲۶۹۵۶۵۹۰	۲۶۹۲۳۵۰۷	۱,۳۴۲	۰,۴۶۶۶	۰	6×10^6	
۵۴۰۷۲۸۴۹	۵۴۰۵۰۹۷۶	۵۴۰۱۶۴۹۹	۱,۵۴	۰,۵۲۷۴	۰	1×10^7	
۳۰۸۲۴۵۲	۳۰۷۵۸۱۱	۳۰۶۸۷۱۷	۱,۳۹۷	۰,۴۰۳۱	۰	1×10^6	توزیع خوشه‌ای
۱۷۶۱۹۰۴۱	۱۷۶۰۶۱۷۴	۱۷۵۹۰۷۳۳	۱,۱۶۳	۰,۴۱۰۳	۰,۰۱۲	6×10^6	
۲۹۲۵۲۹۴۷	۲۹۲۱۹۳۴۶	۲۹۱۹۱۳۶۷	۱,۰۵۴	۰,۴۶۴۵	۰,۰۲۵	1×10^7	

به انواع توزیع داده‌های ورودی، تغییر محسوسی نکرده است.

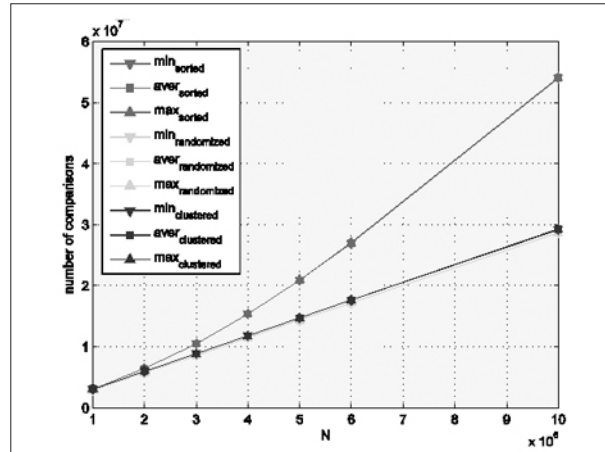
نتیجه جالب و حائز اهمیت در فرایند پیاده‌سازی، گواه بر کارایی الگوریتم با ورودی تصادفی در مقایسه با داده‌های خوشه‌ای و مرتب‌شده است.

۴-۳- مقایسه نتایج نظری و پیاده‌سازی

با توجه به کارایی داده‌های ورودی با توزیع تصادفی، در ادامه نشان می‌دهیم که الگوریتم با ورودی تصادفی در عمل به مراتب رفتار بهتری نسبت به حالت نظری دارد. با تنظیم $\delta = 1$ و $k = \log \log n$ نتایج زیر حاصل شد.

همان‌طور که در شکل (۲) مشاهده می‌کنید تصادفی بودن به‌طور قابل ملاحظه‌ای باعث کاهش تعداد مقایسه‌های انجام‌شده و بهبود خطا شده است. دقت کنید که حتی بیشینه مقدار خطای ورودی تصادفی نیز اختلاف فاحش با مقدار خطای نظری دارد. الگوریتم با ورودی تصادفی در عمل به مراتب رفتار بهتری نسبت به حالت نظری دارد.

۴-۴- تاثیر پارامترهای مختلف بر الگوریتم



شکل ۱: نتایج حاصل از اجرای الگوریتم با ورودی با توزیع تصادفی، مرتب‌شده و خوشه‌ای با تنظیم $\delta = 1$ و $k = \lceil \log \log n \rceil$ و تاثیر آن بر تعداد مقایسه‌های انجام‌شده

توزیع‌های مختلف از داده‌ها را، بهتر مشاهده و تحلیل کرد. برای این کار عملکرد الگوریتم نسبت به ۳ نوع توزیع از داده‌های ورودی مورد بررسی قرار گرفته شده است.

توزیع تصادفی: به تعداد اندازه آرایه‌های ورودی الگوریتم، اعداد به صورت تصادفی طوری تولید می‌شوند که علاوه بر این که در بازه مفروض قرار می‌گیرند، تکراری نیز نباشند.

توزیع مرتب شده: بعد از تولید آرایه‌های تصادفی و بررسی رفتار الگوریتم نسبت به توزیع تصادفی داده‌ها، همان آرایه‌ها مرتب خواهند شد و در اختیار الگوریتم برای بررسی مجدد تغییر رفتار آن، نسبت به این توزیع جدید مرتب شده قرار می‌گیرند.

توزیع خوشه‌ای: برای تولید چنین داده‌هایی، ابتدا \sqrt{n} خوشه با هسته تصادفی ایجاد می‌شود. از هر خوشه نیز به تصادف \sqrt{n} عدد غیرتکراری انتخاب می‌شود. در انتها اجتماع این اعداد، آرایه‌ای از اندازه n با داده‌های \sqrt{n} خوشه‌ای ایجاد می‌کند.

با توجه به شکل (۱)، داده‌ها با توزیع تصادفی نسبت به سایر توزیع‌ها تعداد مقایسه‌ها کمتری را تولید کرده است. پس از این به بعد برای تحلیل رفتار الگوریتم، از ورودی با توزیع تصادفی استفاده خواهیم کرد.

همان‌طور که در جدول (۲) مشاهده می‌کنید خطا نسبت

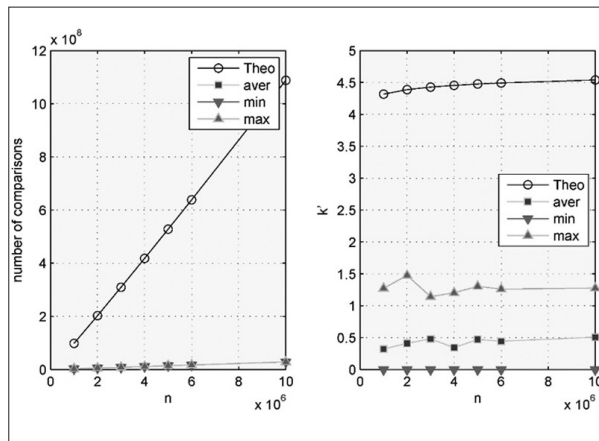
و برای سایر مقادیر k نیز با افزایش میزان k شاهد افت میزان افزایش تعداد مقایسه‌ها هستیم. پس می‌توان نتیجه گرفت؛ هرچقدر میزان k افزایش پیدا می‌کند، تعداد مقایسه‌های انجام‌شده توسط الگوریتم کاهش پیدا می‌کند. اما آیا این میزان کاهش برای همه مقادیر k یکسان است؟ نکته جالب توجه زمانی است که الگوریتم به $k = \lceil \log \log n \rceil$ می‌رسد، دیگر از کاهش چشم‌گیر میزان تعداد مقایسه‌ها خبری نیست.

با نگاهی دقیق‌تر به شکل (۳) می‌توان دریافت که، با وجود خطی بودن همه نمودارها، در شیب، تفاوت وجود دارد. شیب نمودار مربوط به $k = 3$ به مراتب از سایر مقادیر k بیشتر است. روند کاهش شیب تا رسیدن به $k = \lceil \log \log n \rceil$ ادامه دارد. اما بعد از آن روند کاهش شیب نمودارها کند می‌شود. از نگاهی دیگر می‌توان گفت افزایش میزان k بعد از $k = \lceil \log \log n \rceil$ باعث کاهش قابل ملاحظه‌ای در تعداد مقایسه‌ها نمی‌شود. پس در آخر نیز می‌توان نتیجه گرفت؛ $k = \lceil \log \log n \rceil$ یک کران برای کاهش تعداد مقایسه‌ها است. این نتیجه تأییدی بر نتیجه ۲-۱-۳ است.

لازم به ذکر است، کلیه نتایج این قسمت برای تمامی مقادیر δ بررسی شده است و همگی آن‌ها برقرار هستند. فقط برای ترسیم نمودار به صورت دلخواه از $\delta = 1$ استفاده شده است.

تأثیر پارامتر δ

مقدار δ ، آستانه تفاوت است. زمانی که میزان اختلاف دو عدد در حدود δ باشد مقایسه به صورت نادقیق انجام می‌شود، در غیر این صورت مقایسه به صورت دقیق انجام خواهد شد. مقدار δ نیز از ۰,۰۱، ۰,۱، ۰,۲۵، ۰,۵، ۰,۷۵، ۱ و ۱ تغییر می‌کند. نتیجه انجام یک مقایسه نادقیق روی دو عنصر نزدیک به هم، انتخاب یکی از آن دو عنصر به تصادف است، چراکه قدرت تشخیص عنصر بزرگ‌تر را نداریم. در پیاده‌سازی این فرایند با تولید یک عدد طبیعی به تصادف و محاسبه باقی‌مانده آن به ۲ آغاز می‌شود.



شکل ۲: مقایسه نتایج نظری با نتایج حاصل از اجرای الگوریتم با توزیع تصادفی، با تنظیم $\delta = 1$ و $k = \lceil \log \log n \rceil$

در این بخش به بررسی تأثیر پارامترهای k ، δ ، n و s بر نحوه تغییر رفتار الگوریتم در «تعداد مقایسه‌های انجام‌شده» و «خطای خروجی» می‌پردازیم. اندازه مجموعه ورودی می‌باشد که مقادیر ۱، ۲، ۳، ۴، ۵، ۶ و ۱۰ میلیون است.

لازم به ذکر است تغییر اندازه مجموعه ورودی X از ۲۰ میلیون داده به ۲۰۰ میلیون داده، فقط موجب تنوع بیشتر الگوهای ورودی و به تبع آن افزایش احتمال مشاهده الگوهایی با خطای بیشتر از خطای میانگین شد.

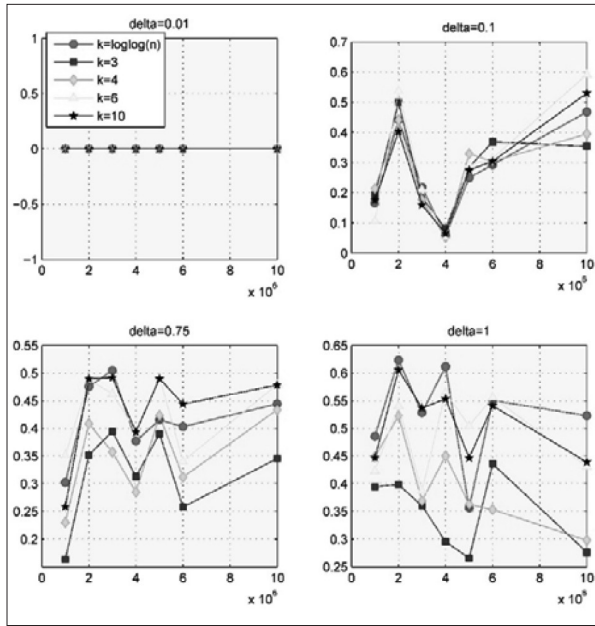
تأثیر پارامتر k

k خطای حاصل از نتایج نظری با مقدار $\log \log n$ است.

این پارامتر در افراز کردن آرایه ورودی نیز شرکت دارد. نحوه تغییر رفتار الگوریتم برای مقادیر $K=3$ (کوچک‌ترین مقدار برای شروع الگوریتم)، ۴، $k = \lceil \log \log n \rceil$ (که به ازای تمام مقادیر ورودی ۵ است)، ۶ و ۱۰ بررسی شد.

نتیجه شهودی قابل استنباط از شکل (۳)، گواه بر خطی بودن تعداد مقایسه‌های انجام‌شده توسط الگوریتم به ازای مقادیر مختلف k است. این نتیجه قابل انتظار، تأییدی بر نتیجه ۲-۱-۱ می‌باشد.

نتیجه دوم با نگاه کلی به رفتار تمام مقادیر k حاصل می‌شود. همان‌طور که در شکل (۳) مشاهده می‌کنید نمودار مربوط به $k = 3$ بیشترین تعداد مقایسه بعد از آن $k = 4$



شکل ۴: تاثیر پارامتر δ بر خطا یا همان K'

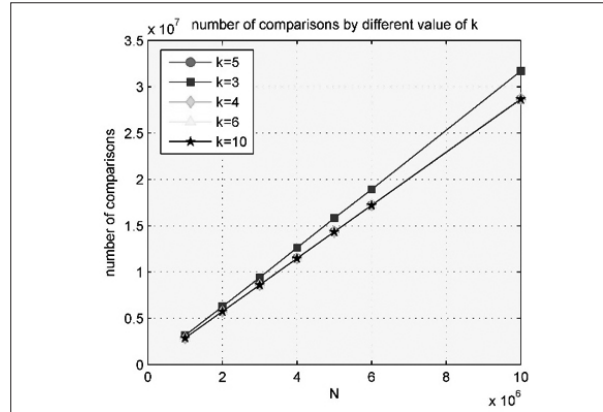
ابتدای کار، دیگر شاهد تغییر رفتار مشابه در میزان افزایش خطا به ازای مقادیر مختلف k نخواهیم بود. با دنبال کردن سیر افزایشی میزان δ ، می‌توان رفتار رو به رشد میزان پراکندگی افزایش خطا را دید.

تاثیر پارامتر s

الگوریتم هر بار s عنصر را به دلخواه انتخاب می‌کند و روی آن‌ها تورنمنت با گردش نوبت برگزار می‌کند. انتخاب‌های متفاوت این s عنصر موجب نتایج متفاوت در خروجی می‌شود. پس اگر نظم یا الگوی خاصی را برای انتخاب این عناصر به الگوریتم ابلاغ کنیم، رفتار الگوریتم به ازای همان چینش خاص مشاهده می‌شود و حتماً روی نتیجه خروجی نیز تاثیر گذار خواهد بود. بنابراین در پیاده‌سازی این s عنصر، به تصادف انتخاب می‌شوند.

در نتایج نظری به صورت پیش فرض مقدار s در الگوریتم 2-MaxFind، $s = \lfloor \sqrt{n} \rfloor$ و در الگوریتم 1-Cover مقدار $s = \lfloor 2\sqrt{n} \rfloor$ در نظر گرفته شده بود. در پیاده‌سازی، مقدار s به ترتیب برای هر کدام از توابع به $s = \lceil \log n \rceil$ و $s = \lceil 2 \log n \rceil$ تغییر دادیم. نتایج این تغییرات در شکل (۵) قابل مشاهده است.

همان‌طور که انتظار داشتیم کاهش اندازه s موجب



شکل ۳: تاثیر پارامتر k بر تعداد مقایسه‌های انجام شده توسط الگوریتم

باقی‌مانده یک عدد به ۲ یا ۰ است یا ۱. اگر باقی‌مانده ۰ بود عدد اول به عنوان برنده مقایسه در نظر گرفته می‌شود و به عنوان خروجی تابع مقایسه‌کننده در اختیار سایر توابع قرار می‌گیرد، در غیر این صورت، عدد دوم برگردانده خواهد شد. مقایسه‌های دقیق نیز مقایسه به صورت معمول انجام می‌شود.

با در نظر گرفتن این‌که اعداد ورودی دارای سه رقم اعشار هستند، زمانی که $\delta = 0.1$ باشد، حداکثر ۱۰ عدد می‌توانند در حدود δ باشند تا مقایسه نادقیق روی آن‌ها اعمال شود. بار دیگر زمانی که δ افزایش پیدا می‌کند و مقدار آن ۰.۱ می‌شود، این تعداد به ۱۰۰، افزایش می‌یابد. افزایش δ باعث می‌شود تعداد اعدادی که در حدود δ هستند، افزایش یابد. این افزایش در جهت افزایش شکل‌گیری مقایسه‌های نادقیق و احتمال افزایش خطاست. همان‌طور که در شکل (۴) مشاهده می‌کنید، در ابتدا با توجه به این‌که δ بسیار کوچک است برای تمامی مقادیر k خطا صفر است. هنگامی که δ افزایش پیدا می‌کند، در تمامی مقادیر k ، این افزایش میزان خطا قابل مشاهده است. نکته جالب توجه در پراکندگی میزان افزایش خطا به ازای k ‌های مختلف است. به نمودار مربوط به $\delta = 0.1$ دقت کنید. رفتار همگن میزان افزایش خطا در k ‌های مختلف، نشان می‌دهد که افزایش خطا به ازای تغییر میزان k بسیار کم است. هرچقدر δ افزایش می‌یابد پراکندگی میزان افزایش خطا بیشتر می‌شود. با افزایش ۱۰۰۰ برابری δ نسبت به

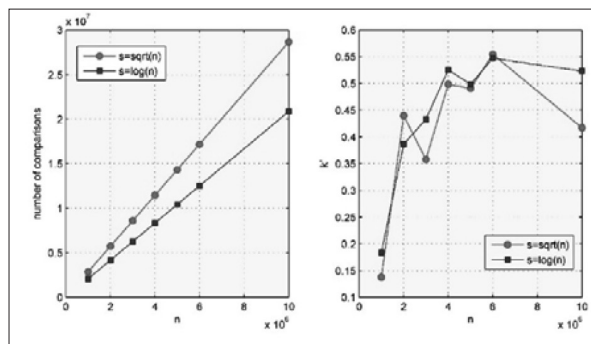
توجه است.

سپاسگزاری

از بهنام ایرانفر و امیر مصری‌خانی جهت راهنمایی در فهم بهتر برخی مفاهیم نظری سپاسگزاری می‌شود.

مراجع

1. Nasibeh Siadaty, Karim Mohammadi, "Design and Simulation of Fault Tolerant Algorithms in Network-on-Chip", Journal of Mathematical Modellings, Vol. 16. No. 2, pp. 33-43, 2010.
2. Yao A C., and Yao F F., "On fault-tolerant networks for sorting", SIAM Journal on Computing, Vol. 14, pp. 120-128. January, 1985.
3. Ravikumar B., Ganesan K, and Lakshmanan K B., "On selecting the largest element in spite of erroneous information", Annual Symposium on Theoretical Aspects of Computer Science, pp. 88-99, Springer, 1987.
4. Ajtai Miklós, Feldman Vitaly, Hassidim Avinatan and Nelson Jelani., "Sorting and selection with imprecise comparisons", ACM Transactions on Algorithms (TALG), Vol. 12, pp. 19, November, 2016.



شکل (۵): تاثیر کاهش مقدار پارامتر S روی خطا و تعداد مقایسه‌ها

کاهش تعداد مقایسه‌ها و افزایش خطا شد.

۵- نتیجه‌گیری

در این مقاله، رفتار الگوریتم محاسبه k -بیشینه در مدل مقایسه نادقیق که ممکن است نتیجه مقایسه داده‌های نزدیک به هم اشتباه باشد، مورد بررسی قرار گرفت. داده‌های مورد استفاده، داده‌های تولید شده تصادفی با سه توزیع مرتب‌شده، یکنواخت و خوشه‌ای است و تعداد مقایسه‌های انجام شده در اجرای الگوریتم و میزان خطای حاصل در بین نمونه‌های مختلف مقایسه گردیدند.

در مجموع، در داده‌های مرتب‌شده، تعداد مقایسه‌های انجام شده به طور ملموسی از سایر توزیع‌ها بیشتر است اما تفاوت ملموسی بین سایر توزیع‌ها در این خصوص مشاهده نشد. همچنین میزان خطا و تعداد مقایسه انجام شده در عمل، به میزان زیادی با نتایج نظری آن‌ها متفاوت است. تعداد مقایسه‌ها برای مقادیر مختلف k نیز بررسی گردید که تفاوت زیادی در تعداد مقایسه‌ها در بین مقادیر مختلف k مشاهده نگردید.

همچنین با تغییر میزان نادقیقی (پارامتر δ) نیز نشان داده شد که به ازای δ های کوچک، میزان خطا تقریباً صفر است و با افزایش میزان δ ، میزان خطا افزایش پیدا می‌کند و مقدار k نیز در میزان خطا تاثیرگذار است.

به‌عنوان کارهای آتی، بررسی الگوریتم‌های پیچیده‌تر نظیر مرتب‌سازی داده‌ها در این مدل در نظر جالب