

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

بررسی الگوریتم‌های موجود برای مسئله کوتاه‌ترین ابرتوالی مشترک و ارائه
الگوریتمی بهبود یافته برای آن

پایان‌نامه کارشناسی ارشد مهندسی کامپیوتر - هوش مصنوعی

فاطمه بحری

استاد راهنما

دکتر سید رسول موسوی



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

پایان‌نامه‌ی کارشناسی ارشد رشته‌ی مهندسی کامپیوتر- هوش مصنوعی خانم فاطمه بحری
تحت عنوان

بررسی الگوریتم‌های موجود برای مسئله کوتاه‌ترین ابرتوالی مشترک و ارائه الگوریتمی بهبود
یافته برای آن

در تاریخ ۱۳۸۹/۱۲/۲۵ توسط کمیته‌ی تخصصی زیر مورد بررسی و تصویب نهایی قرار گرفت.

دکتر سید رسول موسوی

۱- استاد راهنمای پایان‌نامه

دکتر مازیار پالهنگ

۲- استاد داور

دکتر عبدالرضا میرزائی

۳- استاد داور

دکتر سید محمود مدرس هاشمی

سرپرست تحصیلات تکمیلی دانشکده

نخستین سپاس به پیشگاه حضرت دوست، که هر چه داریم از اوست. صمیمانه‌ترین سپاس‌ها را تقدیم می‌کنم به همسر مهربانم به خاطر همدلی‌ها، همراهی‌ها و کمک‌های بی‌دریغش، به پدر و مادر عزیزم که تا ابد وام‌دار محبتشان هستم و به برادران بسیار عزیزم، علیرضا و امیرمحمد. از استاد راهنمای بزرگوارم جناب آقای دکتر موسوی که همواره در مراحل انجام پژوهش از راهنمایی‌های ارزنده ایشان برخوردار بوده‌ام، کمال تشکر و سپاس را دارم.

از اساتید گران‌قدری که در دوران تحصیل افتخار شاگردی‌شان را داشته‌ام، به‌ویژه جناب آقای دکتر صدوری متشکرم. از همه دوستانم که در طول مقطع کارشناسی ارشد در محضرشان بودم و این مدت از زندگیم را به دورانی دوست‌داشتنی تبدیل کردند، به‌ویژه خانم‌ها مهندس عموزادی، پاکیزه، احمدی، طباطبای، رزاقی، ایمانی، بابائی و رضائی کمال تشکر را دارم. و در پایان از همدلی دوستان همیشه همراهم فاطمی، ملیحه، شکوه، فرزانه، نیلوفر، مهناز، مطهره و بهناز سپاسگزارم.

کلیه‌ی حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوری‌های ناشی از تحقیق موضوع این
پایان‌نامه (رساله) متعلق به دانشگاه صنعتی اصفهان
است.

تقدیم به پدر و مادر عزیزتر از جانم
به پاس مهر بی دریغشان

و تقدیم به

مهربانم، همراهم، عزیزم

همسرم

محمد

فهرست مطالب

<u>صفحه</u>	<u>عنوان</u>
هشت	فهرست مطالب.....
۱	چکیده.....
	فصل اول: مقدمه
۲	۱-۱ مقدمه.....
۳	۲-۱ تعریف مسئله SCS.....
۴	۳-۱ کاربردهای مسئله SCS.....
۶	۴-۱ ساختار پایاننامه.....
	فصل دوم: مروری بر الگوریتم های موجود برای مسئله SCS
۷	۱-۲ مقدمه.....
۸	۲-۲ الگوریتم های دقیق.....
۹	۱-۲-۲ الگوریتم برنامه نویسی پویا برای حل مسئله SCS با دو رشته ورودی.....
۱۰	۲-۲-۲ الگوریتم برنامه نویسی پویا برای حل مسئله SCS با n رشته ورودی.....
۱۲	۳-۲ الگوریتم های مکاشفه ای ساده.....
۱۳	۱-۳-۲ الگوریتم ادغام اکثریت.....
۱۳	۲-۳-۲ الگوریتم ادغام اکثریت وزن دار.....
۱۴	۳-۳-۲ الگوریتم ادغام اکثریت آینده نگر.....
۱۵	۴-۳-۲ الگوریتم های حریصانه و مسابقه.....
۱۶	۵-۳-۲ الگوریتم PBS.....
۱۷	۴-۲ الگوریتم های مکاشفه ای مخصوص رشته های DNA.....
۱۷	۱-۴-۲ الگوریتم کمترین ارتفاع و الگوریتم مجموع ارتفاعات.....
۱۸	۲-۴-۲ الگوریتم پس پردازش آینده نگر.....
۱۹	۳-۴-۲ الگوریتم هابل - موریس - وینکلر.....
۱۹	۵-۲ الگوریتم های فرامکاشفه ای.....
۱۹	۱-۵-۲ الگوریتم ژنتیک بر مبنای مکاشفه.....
۲۰	۲-۵-۲ الگوریتم ژنتیک چند هدفه.....
۲۰	۳-۵-۲ الگوریتم ممتیک.....
۲۱	۴-۵-۲ الگوریتم POEMS.....
۲۳	۵-۵-۲ الگوریتم گروه مورچه ها.....
۲۳	۶-۵-۲ روش های تکاملی جریمه کننده و اصلاح کننده.....
۲۴	۷-۵-۲ الگوریتم ترکیبی MA-BS.....
۲۶	۶-۲ الگوریتم های تقریبی.....
۲۶	۱-۶-۲ الگوریتم الفبا.....
۲۶	۲-۶-۲ الگوریتم تقریبی A*.....
۲۷	۳-۶-۲ الگوریتم کاهش - بسط.....
۳۰	۴-۶-۲ الگوریتم DR.....
۳۱	۷-۲ نتیجه گیری..... هشت
	فصل سوم: الگوریتم دقیق A* پیشنهادی

۳۲	۱-۳	مقدمه
۳۳	۲-۳	الگوریتم جستجوی سطحی
۳۴	۱-۲-۳	استفاده از راهبرد شاخه و حد در الگوریتم جستجوی سطحی
۳۶	۲-۲-۳	استفاده از هرس غالبانه در الگوریتم جستجوی سطحی
۳۷	۳-۲-۳	استفاده از راهبرد شاخه و حد و هرس غالبانه در الگوریتم جستجوی سطحی
۳۸	۳-۳	الگوریتم جستجوی عمقی
۳۹	۱-۳-۳	استفاده از راهبرد شاخه و حد در الگوریتم جستجوی عمقی
۴۰	۲-۳-۳	استفاده از هرس غالبانه در الگوریتم جستجوی عمقی
۴۱	۳-۳-۳	استفاده از راهبرد شاخه و حد و هرس غالبانه در الگوریتم جستجوی عمقی
۴۲	۴-۳	الگوریتم A*
۴۳	۱-۴-۳	توابع تخمین مکاشفه ای استفاده شده
۴۶	۲-۴-۳	استفاده از راهبرد شاخه و حد در الگوریتم A*
۴۷	۳-۴-۳	استفاده از هرس غالبانه در الگوریتم A*
۴۸	۴-۴-۳	استفاده از راهبرد شاخه و حد و هرس غالبانه در الگوریتم A*
۴۹	۵-۳	نتیجه گیری
		فصل چهارم: الگوریتم مکاشفه ای پیشنهادی و رویه پس پردازش ارائه شده
۵۰	۱-۴	مقدمه
۵۱	۲-۴	الگوریتم IBS-SCS
۵۲	۱-۲-۴	شرح قدم به قدم الگوریتم IBS-SCS
۵۳	۲-۲-۴	مکاشفه پیشنهادی
۵۶	۳-۲-۴	بررسی پیچیدگی زمانی الگوریتم
۵۷	۳-۴	رویه پس پردازش
۵۷	۱-۳-۴	رویه کاهش الگوریتم DR
۵۹	۲-۳-۴	رویه پس پردازش پیشنهادی
۶۰	۳-۳-۴	تفاوت های رویه پس پردازش پیشنهادی و رویه کاهش DR
۶۰	۴-۴	نتیجه گیری
		فصل پنجم: ارزیابی نتایج
۶۲	۱-۵	مقدمه
۶۳	۲-۵	نتایج الگوریتم A* ارائه شده
۶۴	۱-۲-۵	مقایسه نسخه های مختلف الگوریتم جستجوی سطحی
۶۶	۲-۲-۵	مقایسه نسخه های مختلف الگوریتم جستجوی عمقی
۶۷	۳-۲-۵	مقایسه نسخه های مختلف ارائه شده الگوریتم A*
۷۲	۴-۲-۵	مقایسه الگوریتم A* پیشنهادی با الگوریتم های جستجوی سطحی و جستجوی عمقی
۷۵	۵-۲-۵	مقایسه الگوریتم A* پیشنهادی با الگوریتم برنامه نویسی پویا
۷۵	۳-۵	نتایج الگوریتم IBS-SCS
۷۶	۱-۳-۵	مقایسه الگوریتم IBS-SCS با الگوریتم DR. نه
۸۲	۲-۳-۵	مقایسه الگوریتم IBS-SCS با الگوریتم های PBS و MA-BS
۸۷	۴-۵	نتایج رویه پس پردازش پیشنهاد شده
۸۸	۱-۴-۵	مقایسه نتایج بر روی مجموعه داده تصادفی
۸۹	۲-۴-۵	مقایسه نتایج بر روی مجموعه داده واقعی
۹۰	۵-۵	نتیجه گیری

فصل ششم: نتیجه‌گیری و پیشنهادات

۹۱.....	نتیجه‌گیری.....	۱-۶
۹۲.....	راهکارهای آینده.....	۲-۶
۹۱.....	مراجع.....	

چکیده

مسئله کوتاه‌ترین ابرتوالی مشترک (SCS) از مسائل مهم بهینه‌سازی ترکیبیاتی است. مجموعه‌ای از رشته‌ها به عنوان ورودی به مسئله SCS داده می‌شود و این مسئله به دنبال رشته ای می‌گردد که ابرتوالی تمام رشته‌ها بوده و دارای کوتاه‌ترین طول ممکن باشد. ابرتوالی یک رشته، از درج صفر یا تعداد بیشتری نماد در هر مکان از رشته مورد نظر به دست می‌آید. این مسئله در زمینه‌های مختلفی از جمله فشرده‌سازی داده، بهینه‌سازی پرسش در پایگاه داده‌ها، طرح‌ریزی و زیست‌اطلاعات کاربرد دارد. یکی از کاربردهای مهم این مسئله در زیست‌اطلاعات، کاهش زمان، هزینه و خطا در ساخت یکی از انواع ریزآرایه‌های DNA می‌باشد. ریزآرایه‌های DNA ابزارهای مفیدی هستند که در دسته بندی ژن‌ها، تشخیص ژن‌ها، تصویر سازی بیان ژنی و آشکارسازی SNP، کاربرد دارند.

تا به حال الگوریتم‌های دقیق برنامه‌نویسی پویا و شاخه‌وحد برای حل مسئله SCS ارائه شده‌اند، اما این مسئله هنگامی که تعداد رشته‌ها بیشتر از ۲ باشد یک مسئله NP-Hard می‌باشد، بنابراین تاکنون راه حل بهینه ای که در زمان چندجمله‌ای قابل اجرا باشد برای آن ارائه نشده است و در عمل الگوریتم‌های دقیق برای نمونه‌های بزرگ این مسئله قابل اجرا نیستند. به همین دلیل، اکثر الگوریتم‌های موجود برای حل این مسئله غیر دقیق می‌باشند، بدین معنی که این الگوریتم‌ها از شرط بهینگی پاسخ چشم‌پوشی کرده و سعی می‌کنند که جوابی نزدیک به بهینه در زمانی معقول به دست آورند. تا به حال الگوریتم‌های غیردقیق مختلفی از انواع الگوریتم‌های تقریبی، مکاشفه‌ای ساده و فرا مکاشفه‌ای برای مسئله SCS ارائه شده‌اند. الگوریتم‌های تقریبی الگوریتم‌هایی هستند که جواب حاصل از آنها نسبت به جواب بهینه از حدی که با نسبت تقریب مشخص می‌شود، بدتر نمی‌باشد. الگوریتم‌های مکاشفه‌ای نیز، به الگوریتم‌هایی گفته می‌شود که بر مبنای یک مکاشفه عمل می‌کنند. این الگوریتم‌ها به دو دسته مکاشفه‌ای ساده و فرامکاشفه‌ای تقسیم می‌شوند. برای حل مسئله کوتاه‌ترین ابرتوالی مشترک تا به حال از الگوریتم‌های مکاشفه‌ای ساده و فرامکاشفه‌ای مختلفی استفاده شده است. تعدادی از این الگوریتم‌های مکاشفه‌ای به طور خاص برای رشته‌های DNA طراحی شده‌اند.

در این پایان نامه یک الگوریتم دقیق A^* ارائه شده و برای آن توابع مکاشفه‌ای مختلفی بررسی شده است و در نهایت یک تابع مکاشفه‌ای ترکیبی ارائه شده که با تخمین حدی دقیق‌تر، از سایر توابع مکاشفه‌ای بهتر عمل می‌کند. در این الگوریتم برای کاهش فضای جستجو و در نتیجه کاهش زمان اجرا از راهبرد شاخه و حد و همچنین ایده هرس غالبانه استفاده شده است. با مقایسه این الگوریتم با دیگر الگوریتم‌های دقیق موجود، مشخص شد که الگوریتم ارائه شده منجر به کاهش قابل ملاحظه‌ای در فضای جستجو شده و دارای زمان اجرای بهتری می‌باشد.

همچنین یک الگوریتم مکاشفه‌ای به همراه یک رویه پس پردازش در این پایان نامه ارائه شده است. این الگوریتم مکاشفه‌ای یک الگوریتم جستجوی پرتوی محلی است که مکاشفه‌ای مبتنی بر روابط بازگشتی را به کار گرفته و با استفاده از برنامه‌نویسی پویا یک آرایه از مقادیر احتمالاتی به صورت برنامه‌نویسی پویا تولید می‌کند که برای محاسبه سریع‌تر مقادیر مکاشفه‌ای مورد استفاده قرار می‌گیرد. همچنین این الگوریتم با استفاده از ایده هرس غالبانه فضای جستجو را هرس می‌کند. این الگوریتم با سه الگوریتم برتر اخیر بر روی رشته‌های تصادفی و واقعی مقایسه شده و نتایجی بهتر هم از نظر زمان اجرا و هم طول جواب به دست آورده است. رویه پس پردازش ارائه شده را نیز می‌توان بر روی هر پاسخ غیر بهینه از یک نمونه مسئله SCS، اجرا کرد. این رویه با کاهش طول جواب موجود در حلقه‌ای تکرار شونده، جواب کوتاه‌تری که همچنان ابرتوالی مشترک رشته‌های ورودی می‌باشد، ارائه می‌کند. نتایج به دست آمده مشخص می‌کند که این رویه پس‌پردازش نسبت به رویه پس‌پردازش موجود، در زمان اجرای یکسان، دارای کیفیت جواب بهتری می‌باشد.

کلمات کلیدی: ۱- کوتاه‌ترین ابرتوالی مشترک ۲- الگوریتم A^* ۳- الگوریتم‌های مکاشفه‌ای ۴- جستجوی

پرتوی محلی

فصل اول

مقدمه

۱-۱ مقدمه

مسئله کوتاه ترین ابر توالی مشترک^۱ (SCS) از جمله مسائل مهم بهینه سازی ترکیباتی^۲ است. این مسئله بایستی کوتاه ترین رشته ای را بیابد که ابر توالی هر کدام از رشته های ورودی باشد. ابر توالی رشته x رشته ای است که با اضافه کردن صفر یا تعداد بیشتری نماد، در هر مکان دلخواه از رشته x به دست می آید. برای مثال کوتاه ترین ابر توالی مشترک سه رشته $CGATG$ ، $TCGCA$ و $GATTC$ ، رشته $TCGCATTGC$ می باشد.

مسئله SCS برای ۲ رشته ورودی در زمان چند جمله ای به صورت بهینه [۱] حل می شود، اما هنگامی که تعداد رشته های ورودی بیشتر از ۲ باشد، این مسئله NP-Hard می باشد [۲-۴]. این مسئله حتی در حالت های محدود شده ای که اندازه الفبا ۲ می باشد و یا طول همه رشته ها ۲ می باشد، نیز NP-Hard باقی می ماند [۲-۷]. تا به حال الگوریتم های بهینه ای که برای مسائل NP-Hard ارائه شده اند، بر حسب اندازه ورودی مسئله، دارای زمان اجرایی نمایی می باشند [۸] که این موضوع در مورد مسئله SCS نیز صادق است. الگوریتم های دقیقی که تا به حال برای مسئله SCS ارائه شده اند، الگوریتم برنامه نویسی پویا [۴، ۹] و الگوریتم های شاخه و حد، می باشند. برای حل این مسئله، الگوریتم های غیر دقیق متعددی از انواع تقریبی، مکاشفه ای ساده و فرامکاشفه ای ارائه شده است که در فصل دوم به تفصیل، مورد مطالعه قرار می گیرند.

^۱ Shortest Common Supersequence

^۲ Combinatorial optimization

۲-۱ تعریف مسئله SCS

اگر s یک رشته باشد، طول آن با $|s| = m$ نمایش داده می شود. $s(k)$ برای نمایش k امین نماد رشته s استفاده می شود، که k مقداری بین ۱ و m می باشد. رشته ای که از نماد j ام تا k ام رشته s به دست می آید به طوری که $j \leq k$ ، با نماد $s(j:k)$ نمایش داده می شود، بنابراین برای نمایش رشته شامل k نماد ابتدایی s ، از نماد $s(1:k)$ استفاده می شود. برای نمایش الحاق دو رشته از نماد + استفاده می شود، برای مثال رشته x از الحاق دو رشته x_1 و x_2 به دست می آید: $x = x_1 + x_2$.

فرض کنیم s_1 و s_2 دو رشته باشند و $A_1 = \{i | i \in \mathbb{N}, i \leq |s_1|\}$ و $A_2 = \{i | i \in \mathbb{N}, i \leq |s_2|\}$ که \mathbb{N} مجموعه اعداد صحیح بزرگتر از صفر است. رشته s_2 ابرتوالی رشته s_1 است یا به بیان نمادین $s_1 < s_2$ ، اگر یک تابع یک به یک همانند g از A_1 به A_2 وجود داشته باشد به طوری که:

$$\begin{aligned} \forall k \in A_1, s_1(k) &= s_2(g(k)) \\ \forall k, k' \in A_1, k < k' &\Rightarrow g(k) < g(k') \end{aligned} \quad 1-1$$

تابع g لزوماً یکتا نیست و احتمال دارد چندین تابع که این شروط را برآورده کند، وجود داشته باشد. هر رشته، ابر توالی رشته تهی - که رشته ای به طول صفر است - محسوب می شود. همچنین در صورتی که رشته s_2 ابرتوالی رشته s_1 باشد، می توان گفت که رشته s_1 زیرتوالی رشته s_2 می باشد.

فرض می شود که x یک رشته و S یک مجموعه غیر تهی از رشته ها می باشد. اگر $s_i < x, \forall s_i \in S$ ، آن گاه x ابرتوالی مجموعه رشته S است یا به بیانی نمادین $S < x$. مسئله کوچکترین ابرتوالی مشترک به این صورت تعریف می شود؛ با توجه به یک مجموعه رشته S داده شده، باید رشته x با طول کمینه یافت شود به طوری که $S < x$. اگر تنها یک رشته به عنوان ورودی داده شود، منظور این است که S دارای فقط یک رشته است. الفبایی که رشته های ورودی بر روی آن تعریف شده اند با Σ نمایش داده می شود؛ فرض می کنیم که $|\Sigma| > 1$. از n برای نمایش تعداد رشته های ورودی استفاده می کنیم، به عبارتی دیگر $n = |S|$. از آنجایی که برای $n = 2$ مسئله SCS در زمان خطی حل می شود، فرض می کنیم که $n > 2$ می باشد.

هر کدام از رشته های ورودی مسئله را با حرف s کوچک نشان می دهیم که با زیرنویس ۱ تا n مشخص شده اند یعنی $S = \{s_1, \dots, s_n\}$. از نماد m_i برای نمایش طول s_i استفاده می شود یعنی $m_i = |s_i|$. فرض می شود که $m_i > 0, i = 1, \dots, n$ یعنی هر رشته ورودی حداقل شامل یک نماد می باشد. نماد m نشان دهنده مقدار بیشینه طول رشته های ورودی است یا به عبارتی $m = \max\{m_i, i = 1, \dots, n\}$. از نماد x (شاید همراه با زیرنویس) برای نمایش یک جواب نامزد^۱ استفاده می شود. جواب نامزد x ، جواب نهایی^۱ نامیده می شود

^۱ Candidate solution

اگر x ابرتوالی تمام رشته های ورودی باشد، در غیر این صورت جواب جزئی^۲ نامیده می شود. یک جواب نامزد نهایی در صورتی بهینه است که هیچ جواب نامزد نهایی دیگری با طولی کوتاه تر وجود نداشته باشد.

اگر x یک جواب نامزد جزئی باشد، $p_i(x)$ نشان دهنده ماکزیمم عدد صحیح ممکن k است به طوری که $x \prec s_i(1:k)$ ، به بیانی دیگر $s_i(1:k)$ نشان دهنده بلندترین زیررشته از ابتدای رشته s_i است که توسط x پوشش داده شده است و $p_i(x)$ نشانگر طول این رشته است. رشته ای که از حذف کردن $p_i(x)$ کارا کتر اول رشته s_i به دست می آید، $r_i(x)$ یا به عبارتی دیگر باقیمانده رشته s_i نامیده می شود. $R(x)$ مجموعه تمام باقیمانده های رشته های ورودی است، یعنی $\{r_i(x), i=1, \dots, n\}$. شکل ۱-۱ مثالی برای مشخص شدن نماد $p_i(x)$ و $r_i(x)$ ارائه می کند.

$x = \text{b c a b}$		
$s_1 = \text{b c a } \underbrace{\text{d c f d c}}_{r_1(x)}$	$s_2 = \text{c a b } \underbrace{\text{a a f d c}}_{r_2(x)}$	$s_3 = \text{b a c } \underbrace{\text{d d f c d}}_{r_3(x)}$
$p_1(x) = 3$	$p_2(x) = 3$	$p_3(x) = 2$

شکل ۱-۱- مثالی از $p_i(x)$ و $r_i(x)$ برای جواب جزئی x و سه رشته ورودی s_1 ، s_2 و s_3

۳-۱ کاربردهای مسئله SCS

مسئله SCS در زمینه های مختلفی از جمله فشرده سازی داده^۳ [۱۰]، بهینه سازی پرسش^۴ در پایگاه داده ها [۱۱، ۱۲]، طرح ریزی^۵ [۱۳] و زیست اطلاعات^۶ [۱۴-۱۸] کاربرد دارد. زیست اطلاعات، شاخه ای از علم است که دو شاخه ی علوم کامپیوتر و زیست شناسی را به یکدیگر متصل می کند. یکی از کاربرد های این مسئله در زمینه زیست اطلاعات، مربوط به تولید نوعی ریزآرایه DNA^۷ می باشد. ریزآرایه های DNA شامل هزاران مکان هستند، که هر کدام از آن ها حاوی تعداد زیادی از یک توالی خاص DNA است و در دسته بندی ژن ها^۸، تصویر سازی بیان ژنی^۹، آشکارسازی SNP^{۱۰}، تشخیص ژن ها^{۱۱} و تشخیص بیماری ها کاربرد دارند [۱۹، ۲۰]. ریزآرایه

¹ Final solution
² Partial solution
³ Data compression
⁴ Query optimization
⁵ Planning
⁶ Bioinformatics
⁷ DNA microarray
⁸ Gene clustering
⁹ Gene expression profiling
¹⁰ Single Nucleotide Polymorphism detection
¹¹ Gene identification

های DNA، به دو نوع ریزآرایه cDNA و ریزآرایه اولیگونوکلوئوتید^۱ تقسیم می شوند. ریزآرایه های اولیگونوکلوئوتید هزینه ساخت بیشتری دارند، اما در آن ها پیوند متقاطع^۲ کمتر رخ می دهد [۲۱] و به همین دلیل دقیق تر می باشند. این نوع ریزآرایه برای تشخیص چندریختی^۳ می تواند به کار رود [۱۹].

ریزآرایه های اولیگونوکلوئوتید معمولاً به روش لیتوگرافی نوری^۴ تولید می شوند [۲۲]. در این روش، در هر قدم، یک نوکلئوتید به بعضی از مکان های ریزآرایه اضافه می شود و به انتهای توالی آن مکان، الحاق می شود. نوکلئوتید واحد سازنده اسیدهای نوکلئیک می باشد و بنا به نوعش با یکی از نمادهای^۵ A،^۶ C،^۷ G یا^۸ T نمایش داده می شود. برای اضافه کردن نوکلئوتید در هر دور، یک نقاب^۹ برای محافظت از مکان هایی از ریزآرایه که نایبستی اشعه فرابنفش را دریافت کنند، استفاده می شود. اشعه فرابنفش در مکان های بدون نقاب نفوذ کرده و اولیگونوکلوئوتید را در آنجا فعال می کند [۱۹]. هزینه ساخت نقاب های مورد استفاده در این روش بالا می باشد و بنابراین کاهش مراحل تولید ریزآرایه، به کاهش هزینه ساخت و البته کاهش زمان مورد نیاز برای ساخت ریزآرایه منجر می شود. به دلیل این که پوشاندن مکان های ریزآرایه با نقاب، بدون خطا نمی باشد، با کاهش مراحل می توان به کاهش خطا در ساخت ریزآرایه نیز کمک کرد. مسئله یافتن کمترین مراحل مورد نیاز برای ساخت ریزآرایه ای از تعدادی توالی DNA را می توان با مسئله SCS مدل کرد، به این صورت که هر توالی DNA به عنوان یک رشته ورودی در نظر گرفته شود. کوتاه ترین ابر توالی مشترک این رشته ها، ترتیب اضافه کردن نوکلئوتیدها را (که منجر به کمترین مراحل ساخت می شود) نمایش می دهد. بنابراین استفاده از مسئله SCS موجب کاهش هزینه، زمان و خطا در ساخت ریزآرایه های اولیگونوکلوئوتید می شود [۱۵-۱۷].

به عنوان مثالی دیگر از کاربرد SCS می توان به طراحی خط تولیدی که در آن قطعات مختلفی تولید می شوند، اشاره کرد. برای تولید یک قطعه لازم است دنباله ای از عمل ها که به نوع آن قطعه بستگی دارد بر روی آن صورت بگیرد. خط تولید نیز بایستی به صورت دنباله ای از دستگاه ها طراحی شود، که هر دستگاه قابلیت انجام فقط یک عمل را بر روی قطعه ها دارد. این مسئله با استفاده از مسئله SCS مدل می شود، به این صورت که دنباله عملیات مورد نیاز برای هر قطعه، یک رشته ورودی در نظر گرفته می شود و کوتاه ترین ابر توالی

¹ Oligonucleotide

² Cross hybridization

³ Polymorphism

⁴ Photolithography

⁵ Adenine

⁶ Cytosine

⁷ Guanine

⁸ Thymine

⁹ Mask

مشترک این رشته‌ها، کوتاه‌ترین دنباله از دستگاه‌ها را مشخص می‌کند که تمام قطعات با این دنباله از دستگاه‌ها می‌توانند تولید شوند.

۴-۱ ساختار پایان‌نامه

در این تحقیق، یک الگوریتم دقیق A^* پیشنهاد شده است و به منظور هرس فضای جستجو و کاهش زمان اجرا، در آن از راهبرد شاخه و حد و نیز هرس غالبانه^۱ استفاده شده است. در این الگوریتم از توابع مکاشفه ای مختلفی استفاده شده است و در نهایت یک تابع مکاشفه ای ترکیبی ارائه شده است که نسبت به بقیه توابع مکاشفه ای نتایج بهتری به دست می‌آورد. الگوریتم A^* مطرح شده، نسبت به سایر الگوریتم‌های دقیق از جمله برنامه نویسی پویا دارای زمان اجرای بهتری می‌باشد.

همچنین در تحقیق ارائه شده، یک الگوریتم غیر دقیق مکاشفه ای به نام IBS-SCS برای مسئله SCS ارائه شده است. الگوریتم IBS-SCS یک الگوریتم جستجوی پرتوی محلی^۲ می‌باشد که از مکاشفه ای جدید مبتنی بر روابط بازگشتی و نیز هرس غالبانه برای کاهش فضای جستجو استفاده می‌کند. این الگوریتم در مقایسه با بهترین الگوریتم‌های غیر دقیق برای SCS (که الگوریتم‌های DR [۲۳]، PBS [۲۴] و MA-BS [۲۵] می‌باشند) نتایج بهتری در مدت زمان اجرای کمتر به دست می‌آورد. همچنین یک رویه پس پردازش پیشنهاد شده است که با اجرای این رویه پس پردازش بر روی جواب به دست آمده از هر الگوریتم غیر دقیق حل SCS، جواب کوتاه تری به دست می‌آید. رویه پس پردازش پیشنهاد شده از مرحله کاهش در الگوریتم DR [۲۳] الهام گرفته شده است و نسبت به آن در زمان اجرای یکسان، نتایج بهتری کسب می‌کند.

ساختار این پایان‌نامه بدین ترتیب است: در فصل دوم مروری بر الگوریتم‌هایی که تا به حال برای مسئله SCS ارائه شده اند، صورت می‌گیرد. سپس در فصل سوم الگوریتم A^* پیشنهادی ارائه می‌گردد. الگوریتم مکاشفه ای IBS-SCS و رویه پس پردازش پیشنهادی در فصل چهارم به تفصیل بیان می‌گردند. در فصل پنجم ابتدا نتایج الگوریتم A^* با الگوریتم‌های دقیق پیاده سازی شده مقایسه شده، سپس نتایج الگوریتم IBS-SCS با بهترین الگوریتم‌های موجود مقایسه شده و در پایان مقایسه ای بین نتایج رویه پس پردازش پیشنهادی و رویه کاهش الگوریتم DR صورت می‌گیرد. در فصل ششم نیز جمع بندی این تحقیق و پیشنهاداتی برای آینده ارائه می‌شود.

^۱ Dominance pruning

^۲ Local Beam Search

فصل دوم

مروری بر الگوریتم‌های موجود برای مسئله SCS

۱-۲ مقدمه

الگوریتم‌های موجود برای مسئله SCS را می‌توان به دو دسته دقیق و غیر دقیق تقسیم کرد، که در این فصل مورد مطالعه قرار می‌گیرند. الگوریتم‌های دقیق، الگوریتم‌هایی هستند که مساله را به صورت دقیق حل کرده و پاسخ بهینه مسئله را می‌یابند. این الگوریتم‌ها بر روی نمونه‌های بزرگ از مسائل NP-Hard، قابل اجرا نمی‌باشند. به همین علت الگوریتم‌های دقیق برای مسئله SCS به ندرت مورد بررسی قرار گرفته‌اند. در بخش ۲-۲، این الگوریتم‌های دقیق مطالعه می‌شوند.

الگوریتم‌های غیردقیق با برداشتن شرط بهینگی، در زمان نسبتاً مناسب، جوابی نزدیک به بهینه به دست می‌آورند. الگوریتم‌های غیردقیق دو دسته هستند. دسته اول الگوریتم‌های تقریبی می‌باشند. این الگوریتم‌ها، تضمین می‌کنند که جواب آنها نسبت به جواب بهینه از حدی که نسبت تقریب^۱ نام دارد، بدتر نیست. الگوریتم‌های تقریبی ارائه شده برای مسئله SCS در بخش ۲-۶ مورد مطالعه قرار می‌گیرند. دسته دوم، الگوریتم‌های مکاشفه‌ای نام دارند. این الگوریتم‌ها، هیچ نسبت تقریبی را تضمین نمی‌کنند، اما معمولاً جواب‌های خوبی در زمانی مناسب به دست می‌آورند. الگوریتم‌های مکاشفه‌ای به دو دسته مکاشفه‌ای ساده و فرا مکاشفه‌ای تقسیم می‌شوند. الگوریتم‌های مکاشفه‌ای ساده معمولاً از یک تابع حریصانه^۲ استفاده می‌کنند و با نام حریصانه نیز شناخته

^۱ Approximation ratio

^۲ Greedy

می‌شوند. تا به حال تعدادی الگوریتم مکاشفه‌ای ساده برای مسئله SCS ارائه شده است، که از آن جمله می‌توان به الگوریتم‌های مکاشفه‌ای ساده‌ای که برای رشته‌های DNA طراحی شده‌اند، اشاره کرد. این الگوریتم‌های مکاشفه‌ای ساده و مکاشفه‌ای مخصوص رشته‌های DNA به ترتیب در بخش‌های ۲-۳ و ۲-۴، مورد بررسی قرار گرفته‌اند. الگوریتم‌های فرامکاشفه‌ای روش‌هایی محاسباتی سطح بالایی هستند که به صورت تکراری سعی در بهبود جواب مسئله دارند، ولی یافتن جواب بهینه را تضمین نمی‌کنند. این الگوریتم‌ها در بخش ۲-۵ تشریح شده‌اند.

علاوه بر نسخه اصلی مسئله SCS، تا به حال نسخه‌های خاص از مسئله SCS نیز مورد توجه قرار گرفته‌اند. در [۲۶] برای دو نسخه خاص از مسئله SCS، الگوریتم‌های تقریبی ارائه شده است. در نسخه اول، فرض شده است که تمام رشته‌ها دارای طول ۲ می‌باشند. در نسخه دوم نه تنها فرض شده است همه رشته‌ها دارای طول ۲ می‌باشند، بلکه هر نماد از الفبا در S (مجموعه رشته‌های ورودی) حداکثر ۳ بار تکرار شده است و برای آن یک الگوریتم با نسبت تقریب ۷/۶ ارائه شده است. در [۲۷] نیز حالت محدود شده مسئله SCS که همه رشته‌ها دارای طول ۲ می‌باشند، مورد بررسی قرار گرفته است.

۲-۲ الگوریتم‌های دقیق

در [۴] الگوریتم برنامه‌نویسی پویا برای حل مسئله SCS بیان شده است. همچنین در [۹] به بررسی راه‌حل‌های دقیق مسئله SCS با استفاده از الگوریتم برنامه‌نویسی پویا و الگوریتم‌های شاخه و حد پرداخته شده است. الگوریتم‌های برنامه‌نویسی پویا فقط بر روی تعداد کمی رشته داده شده موفق عمل می‌کنند و در صورت زیاد شدن تعداد رشته‌ها، زمان و فضای حافظه بسیار زیادی نیاز دارند. کوتاه‌ترین ابرتوالی مشترک دو رشته با طول m توسط برنامه‌نویسی پویا با پیچیدگی زمانی و مکانی $o(m^2)$ محاسبه می‌شود. این الگوریتم مسئله‌ای با n رشته با طول m را با پیچیدگی زمانی و مکانی $o(m^n)$ محاسبه می‌کند [۴]. الگوریتم‌های شاخه و حد نیز به زمان زیادی برای محاسبه جواب بهینه نیاز دارند، مگر آن که اندازه الفبا بسیار کوچک باشد [۹]. در این بخش، ابتدا الگوریتم برنامه‌نویسی پویا برای حل مسئله SCS با دو رشته ورودی و سپس این الگوریتم برای حل مسئله SCS با n رشته ورودی بر اساس [۴] تشریح می‌شود.

۱-۲-۲ الگوریتم برنامه نویسی پویا برای حل مسئله SCS با دو رشته ورودی

کوتاه ترین ابر توالی مشترک دو رشته x و y توسط برنامه نویسی پویا، در زمان چند جمله ای به دست می آید. در صورتی که $length(i, j)$ نشان دهنده طول کوتاه ترین ابر توالی مشترک دو رشته $x(1:i)$ و $y(1:j)$ باشد، رابطه بازگشتی ۱-۲ زیر برای یافتن کوتاه ترین ابر توالی مشترک دو رشته x و y به کار می رود [۴].

$$length(i, j) = \begin{cases} j & \text{if } i = 0 \\ i & \text{if } j = 0 \\ length(i-1, j-1) + 1 & \text{if } i, j > 0, x(i) = y(j) \\ \min\{length(i, j-1) + 1, length(i-1, j) + 1\} & \text{if } i, j > 0, x(i) \neq y(j) \end{cases} \quad 1-2$$

شبه کدهای ۱-۲ و ۲-۲ طریقه یافتن کوتاه ترین ابر توالی مشترک دو رشته x و y را به وسیله رابطه

بازگشتی فوق نشان می دهد.

```

find-SCS(x,y)
{
    for(i=0:|x|)
        length(i,0)=i
    for(j=0:|y|)
        length(j,0)=j
    for(i=1:|x|)
        for(j=1:|y|)
            {
                if(x(i)==y(j))
                {
                    length(i,j)=length(i-1,j-1)+1
                    Previous(i,j)=1
                }
                else
                {
                    if(length(i,j-1)<length(i-1,j))
                    {
                        length(i,j)=length(i,j-1)+1
                        Previous(i,j)=2
                    }
                    else
                    {
                        length(i,j)=length(i-1,j)+1
                        Previous(i,j)=3
                    }
                }
            }
        }
    }
back-tracking(|x|,|y|)
}

```

شبه کد ۱-۲- تابع find-SCS برای یافتن SCS دو رشته x و y با برنامه نویسی پویا [۴]

در تابع find-SCS به وسیله برنامه نویسی پویا، طول ابر توالی مشترک x و y در ماتریس $length$ ذخیره می شود. البته ماتریس دیگری به نام $previous$ هم با اندازه ماتریس $length$ در نظر گرفته شده است. در عنصر $previous(i, j)$ مشخص می شود که برای به دست آوردن $length(i, j)$ ، از کدام یک از عناصر $length(i-1, j)$ و $length(i-1, j-1)$ و یا $length(i, j-1)$ استفاده شده است. در پایان که تمام عناصر ماتریس $length$ محاسبه شد، تابع $back - tracking$ با پارامترهای $|x|$ و $|y|$ فرا خوانده می شود.

back-tracking(i,j)

```

{
    if(i==0&j==0)
        return SCS;
    If(i==0)
    {
        SCS=y(0:j)+SCS
        return SCS
    }
}

```

شبهه کد ۲-۲- تابع *back-tracking* برای به دست آوردن *SCS* با استفاده از ماتریس *previous* [۴]

```

If(j==0)
{
    SCS=x(0:i)+SCS
    return SCS
}
switch(previous(i,j))
{
    case1:
    {
        SCS=x(i)+SCS // or SCS=y(j)+SCS
        back-tracking(i-1,j-1)
    }
    case2:
    {
        SCS=y(j)+SCS
        back-tracking(i,j-1)
    }
    case3:
    {
        SCS=x(i)+SCS
        back-tracking(i-1,j)
    }
}
}

```

ادامه شبهه کد ۲-۲- تابع *back-tracking* برای به دست آوردن *SCS* با استفاده از ماتریس *previous* [۴]

تابع *back-tracking* یک تابع بازگشتی است که رشته ابرتوالی x و y را با استفاده از مقادیر ماتریس *previous*، به دست می آورد. با توجه به مقدار $previous(i, j)$ ، در فراخوانی $back-tracking(i, j)$ نماد $x(i)$ و یا $y(j)$ به *SCS* اضافه شده و تابع *back-tracking* با پارامترهای جدید فرا خوانی می شود. هنگامی که j برابر با صفر شود، *SCS* ابرتوالی رشته y است، بنابراین قسمت باقیمانده از رشته x به *SCS* الحاق می شود یعنی $SCS = x(0:i) + SCS$ و مقدار *SCS* بازگردانده می شود. همین طور هنگامی که i برابر با صفر شود، *SCS* ابرتوالی رشته x است، بنابراین قسمت باقیمانده از رشته y به *SCS* الحاق می شود یعنی $SCS = y(0:j) + SCS$ و مقدار *SCS* بازگردانده می شود. و در صورتی که i و j هر دو برابر صفر باشند، *SCS* ابرتوالی هر دو رشته x و y می باشد و بازگردانده می شود. در بخش بعدی، یافتن ابرتوالی n رشته با روش برنامه نویسی پویا بیان می شود.

۲-۲-۲ الگوریتم برنامه نویسی پویا برای حل مسئله *SCS* با n رشته ورودی

در [۴] به صورت خلاصه، یافتن ابرتوالی مشترک n رشته با روش برنامه نویسی پویا بیان شده است که در اینجا تشریح می شود. ابتدا بعضی تعاریف و نمادها بیان می شود. فرض می شود که هر یک از رشته های s_i ،