



دانشگاه علم و صنعت ایران

دانشکده مهندسی کامپیوتر

طراحی و تحلیل الگوریتم‌ها

تمرین ۶*

مطهره میرزایی
محمدجواد پیرهادی
سید صالح اعتمادی

نیم‌سال دوم ۹۸-۹۹

@MohammadJavad_Pirhadi @mirzaei2114	تلگرام
fb_A6	نام شاخه
A6	نام پروژه/پوشه/پول ریکوست
۱۳۹۹/۱/۲۳	مهلت تحویل

توضیحات کلی تمرین

۱. ابتدا مانند تمرین های قبل، یک پروژه به نام A6 بسازید.
۲. کلاس هر سوال را به پروژه‌ی خود اضافه کنید و در قسمت مربوطه کد خود را بنویسید. هر کلاس شامل دو متد اصلی است:
 - متد اول: تابع Solve است که شما باید الگوریتم خود را برای حل سوال در این متد پیاده سازی کنید.
 - متد دوم: تابع Process است که مانند تمرین های قبلی در TestCommon پیاده سازی شده است. بنابراین با خیال راحت سوال را حل کنید و نگران تابع Process نباشید! زیرا تمامی پیاده سازی ها برای شما انجام شده است و نیازی نیست که شما کدی برای آن بزنید.
۳. اگر برای حل سوالی نیاز به تابع های کمکی دارید؛ می توانید در کلاس مربوط به همان سوال تابع تان را اضافه کنید.

اکنون که پیاده سازی شما به پایان رسیده است، نوبت به تست برنامه می رسد. مراحل زیر را انجام دهید.

 ۱. یک UnitTest برای پروژه‌ی خود بسازید.
 ۲. فولدر TestData که در ضمیمه همین فایل قرار دارد را به پروژه‌ی تست خود اضافه کنید.
 ۳. فایل GradedTests.cs را به پروژه‌ی تستی که ساخته اید اضافه کنید.

توجه:

برای اینکه تست شما از بهینه سازی کامپایلر دات نت حداکثر بهره را ببرد زمان تست ها را روی بیلد Release امتحان کنید، در غیر اینصورت ممکن است تست های شما در زمان داده شده پاس نشوند.

```

1 using Microsoft.VisualStudio.TestTools.UnitTesting;
2 using A6;
3 using System;
4 using System.Collections.Generic;
5 using System.Linq;
6 using System.Text;
7 using System.Threading.Tasks;
8 using TestCommon;
9
10 namespace A6.Tests
11 {
12     [TestClass()]
13     [DeploymentItem("TestData", "A6_TestData")]
14     public class GradedTests
15     {
16         [TestMethod(), Timeout(500)]
17         public void SolveTest_Q1ConstructBWT()
18         {
19             RunTest(new Q1ConstructBWT("TD1"));
20         }
21
22         [TestMethod(), Timeout(1000)]
23         public void SolveTest_Q2ReconstructStringFromBWT()
24         {
25             RunTest(new Q2ReconstructStringFromBWT("TD2"));
26         }
27
28         [TestMethod(), Timeout(500)]
29         public void SolveTest_Q3MatchingAgainCompressedString()
30         {
31             RunTest(new Q3MatchingAgainCompressedString("TD3"));
32         }
33
34         [TestMethod(), Timeout(1000)]
35         public void SolveTest_Q4ConstructSuffixArray()
36         {
37             RunTest(new Q4ConstructSuffixArray("TD4"));
38         }
39         public static void RunTest(Processor p)
40         {
41             TestTools.RunLocalTest("A6", p.Process, p.TestDataName, p.Verifier,
42                 VerifyResultWithoutOrder: p.VerifyResultWithoutOrder,
43                 excludedTestCases: p.ExcludedTestCases);
44         }
45     }
46 }

```

۱ بدست آوردن تبدیل Burrows Wheeler برای یک رشته^۱

تبدیل Burrows Wheeler یک متن، نمادهای متن را به گونه‌ای تغییر می‌دهد تا به راحتی قابل فشرده‌سازی شود. همچنین، این تبدیل برگشت‌پذیر نیز هست. یعنی از روی تبدیل BW یک متن می‌توانیم متن اصلی را بدست آوریم. کاربرد این تبدیل علاوه بر فشرده سازی متن، در حل کردن الگوریتم تطبیق چندگانه الگوها و ... است.

بدست آوردن تبدیل Burrows Wheeler برای یک رشته به صورت مقابل تعریف می‌شود: اول تمام چرخش‌های متناوب ممکن از متن را تشکیل دهید. یک چرخش متناوب از جداکردن پسوندی از رشته از انتهای رشته و اضافه کردن آن به ابتدای رشته بدست می‌آید. سپس این رشته‌ها را به ترتیب الفبایی مرتب کنید تا یک ماتریس M با سایز $|TEXT| \times |TEXT|$ تشکیل شود. BWT برای رشته ورودی برابر با ستون آخر ماتریس M است.

در این سوال شما باید الگوریتمی بنویسید که BWT برای یک رشته ورودی که با حروف T G C A ساخته شده است و با نماد \$ پایان می‌یابد را برگرداند. سایز رشته ورودی بین ۱ تا ۱۰۰۰ کاراکتر است.

توجه:

برای دیباگ و تست کردن جواب‌های خود می‌توانید از این لینک استفاده کنید.

ورودی نمونه	خروجی نمونه
AA\$	AA\$

$$M(\text{Text}) = \begin{bmatrix} \$ & A & A \\ A & \$ & A \\ A & A & \$ \end{bmatrix}$$

شکل ۱: نمونه اول

^۱Construct the Burrows–Wheeler Transform of a String

ورودی نمونه	خروجی نمونه
ACACACAC\$	CCCC\$AAAA

$$M(\text{Text}) = \begin{bmatrix} \$ & A & C & A & C & A & C & A & C \\ A & C & \$ & A & C & A & C & A & C \\ A & C & A & C & \$ & A & C & A & C \\ A & C & A & C & A & C & \$ & A & C \\ A & C & A & C & A & C & A & C & \$ \\ C & \$ & A & C & A & C & A & C & A \\ C & A & C & \$ & A & C & A & C & A \\ C & A & C & A & C & \$ & A & C & A \\ C & A & C & A & C & A & C & \$ & A \end{bmatrix}$$

شکل ۲: نمونه دوم

ورودی نمونه	خروجی نمونه
AGACATA\$	ATG\$CAAA

$$M(\text{Text}) = \begin{bmatrix} \$ & A & G & A & C & A & T & A \\ A & \$ & A & G & A & C & A & T \\ A & C & A & T & A & \$ & A & G \\ A & G & A & C & A & T & A & \$ \\ A & T & A & \$ & A & G & A & C \\ C & A & T & A & \$ & A & G & A \\ G & A & C & A & T & A & \$ & A \\ T & A & \$ & A & G & A & C & A \end{bmatrix}$$

شکل ۳: نمونه سوم

```

1 using System;
2 using TestCommon;
3
4 namespace A6
5 {
6     public class Q1ConstructBWT : Processor
7     {
8         public Q1ConstructBWT(string testDataName)
9         : base(testDataName) { }
10
11         public override string Process(string inStr) =>
12             TestTools.Process(inStr, (Func<String, String>)Solve);
13
14         /// <summary>
15         /// Construct the Burrows-Wheeler transform of a string
16         /// </summary>
17         /// <param name="text"> A string Text ending with a "$" symbol </param>
18         /// <returns> BWT(Text) </returns>
19         public string Solve(string text)
20         {
21             throw new NotImplementedException();
22         }
23     }
24 }

```

۲ بدست آوردن رشته از روی BWT ۲

در سوال قبل با چگونه بدست آوردن BWT برای یک رشته آشنا شدید و دیدیم این عمل برای فشرده‌سازی متون انجام می‌شود. برای تکمیل فرآیند فشرده‌سازی لازم است تا بتوانیم برعکس این فرآیند را نیز انجام دهیم. یعنی بتوانیم از روی BWT یک رشته، رشته اصلی را بدست آوریم. در این سوال باید الگوریتمی بنویسید که رشته اصلی را از روی BWT رشته بدست آورد. ورودی الگوریتم یک رشته شامل یک نماد \$ است که با حروف G T C A ساخته شده است و سایز رشته ورودی بین ۱ تا ۱۰۰۰۰۰۰ کاراکتر است. خروجی الگوریتم آن رشته اصلی است که تبدیل BW آن در ورودی داده شده است.

ورودی نمونه	خروجی نمونه
AC\$A	ACA\$

$$M(\text{Text}) = \begin{bmatrix} \$ & A & C & A \\ A & \$ & A & C \\ A & C & A & \$ \\ C & A & \$ & A \end{bmatrix}$$

شکل ۴: نمونه اول

ورودی نمونه	خروجی نمونه
AGGGAA\$	GAGAGA\$

$$M(\text{Text}) = \begin{bmatrix} \$ & G & A & G & A & G & A \\ A & \$ & G & A & G & A & G \\ A & G & A & \$ & G & A & G \\ A & G & A & G & A & \$ & G \\ G & A & \$ & G & A & G & A \\ G & A & G & A & \$ & G & A \\ G & A & G & A & G & A & \$ \end{bmatrix}$$

شکل ۵: نمونه دوم

```

1 using System;
2 using TestCommon;
3
4 namespace A6
5 {
6     public class Q2ReconstructStringFromBWT : Processor
7     {
8         public Q2ReconstructStringFromBWT(string testDataName)
9         : base(testDataName) { }
10
11         public override string Process(string inStr) =>
12         TestTools.Process(inStr, (Func<String, String>)Solve);
13
14         /// <summary>
15         /// Reconstruct a string from its Burrows-Wheeler transform
16         /// </summary>
17         /// <param name="bwt"> A string Transform with a single "$" sign </param>
18         /// <returns> The string Text such that BWT(Text) = Transform.
19         /// (There exists a unique such string.) </returns>
20         public string Solve(string bwt)
21         {
22             throw new NotImplementedException();
23         }
24     }
25 }

```


۳ تطبیق الگوها با استفاده از رشته‌های فشرده شده^۳

یکی دیگر از قابلیت‌های جالب BWT، این است که به ما این قابلیت را می‌دهد که بدون decompress کردن رشته بتوانیم مسئله Pattern Matching را حل کنیم. بدین ترتیب حافظه کمتری مصرف می‌شود. الگوریتم BWMatching تعداد تکرار شدن الگوها را در متن می‌شمارد. اطلاعاتی که این الگوریتم در اختیار دارد تنها ستون اول ماتریس، ستون آخر (یعنی BWT) و یک mapping از ابتدا به انتها است. در این سوال شما باید الگوریتمی بنویسید که بتواند الگوها را در یک متن فشرده شده پیدا کند. در خط اول فایل ورودی یک رشته که فرم BWT متن است وجود دارد. این رشته با نمادهای G C T A ساخته شده است و طول آن بین ۱ تا ۱۰۰۰ حرف است و در آن نماد \$ نیز وجود دارد. در خط بعدی یک عدد صحیح بین ۱ تا ۵۰۰۰ وجود دارد که نشان‌دهنده ی تعداد الگوهاست. در n خط بعدی الگوهایی که باید در رشته ورودی یافت شوند وجود دارد. شما باید در فایل خروجی لیستی از اعداد صحیح که هر یک نشان‌دهنده ی تعداد تکرار الگوی متناظرش در متن است را برگردانید.

ورودی نمونه	خروجی نمونه
AGGGAA\$ 1 GA	3

در این مثال متن GAGAGA\$ بوده که الگوی GA سه بار در آن تکرار شده است.

ورودی نمونه	خروجی نمونه
ATT\$AA 2 ATA A	2 3

در این مثال متن ATATA\$ بوده که شامل ۲ الگوی ATA و سه الگوی A است.

ورودی نمونه	خروجی نمونه
AT\$TCTATG 2 TCT TATG	0 0

در این مثال متن ATCGTTTA بوده که شامل هیچ یک از الگوهای داده شده در ورودی نیست.

```

1 using System;
2 using TestCommon;
3
4 namespace A6
5 {
6     public class Q3MatchingAgainCompressedString : Processor
7     {
8         public Q3MatchingAgainCompressedString(string testDataName)
9         : base(testDataName) { }
10
11         public override string Process(string inStr) =>
12         TestTools.Process(inStr, (Func<String, long, String[], long[]>)Solve);
13
14         /// <summary>
15         /// Implement BetterBWMatching algorithm
16         /// </summary>
17         /// <param name="text"> A string BWT(Text) </param>
18         /// <param name="n"> Number of patterns </param>
19         /// <param name="patterns"> Collection of n strings Patterns </param>
20         /// <returns> A list of integers, where the i-th integer corresponds
21         /// to the number of substring matches of the i-th member of Patterns
22         /// in Text. </returns>
23         public long[] Solve(string text, long n, String[] patterns)
24         {
25             throw new NotImplementedException();
26         }
27     }
28 }

```

۴ تشکیل Suffix Array برای یک رشته^۴

در تمرین قبل با Tree ها Suffix آشنا شدیم و دیدیم در عمل حافظه زیادی اشغال می‌کنند. آرایه پیشوندی یک راه حل جایگزین برای درخت پیشوندی است که از نظر حافظه بهینه عمل می‌کند. به عنوان مثال برای ذخیره درخت پیشوندی ژنوم انسان نیاز به ۶۰ گیگ حافظه داریم در حالی که آرایه پیشوندی تنها ۱۲ گیگ حافظه اشغال می‌کند. برای ساختن Suffix Array یک رشته به این صورت عمل می‌کنیم که تمامی پسوندهای رشته ورودی را به ترتیب الفبایی (با فرض \$ بعنوان اولین نماد الفبا) مرتب می‌کنیم. آرایه پیشوندی، لیستی از ایندکس اولین نماد از این پسوندهای مرتب شده در رشته ورودی است. در این سوال باید آرایه پیشوندی برای یک رشته ورودی ساخته شده با حروف G T C A که با نماد \$ تمام می‌شود را تشکیل دهید.

توجه:

برای دیباگ و تست کردن جواب‌های خود می‌توانید از این لینک استفاده کنید.

ورودی نمونه	خروجی نمونه
GAC\$	3 1 2 0

Sorted suffixes:

3 \$
1 AC\$
2 C\$
0 GAC\$

شکل ۶: نمونه اول

ورودی نمونه	خروجی نمونه
GAGAGAGA\$	8 7 5 3 1 6 4 2 0

Sorted suffixes:

8 \$
7 A\$
5 AGA\$
3 AGAGA\$
1 AGAGAGA\$
6 GA\$
4 GAGA\$
2 GAGAGA\$
0 GAGAGAGA\$

شکل ۷: نمونه دوم

ورودی نمونه	خروجی نمونه
AACGATAGCGGTAGA\$	15 14 0 1 12 6 4 2 8 13 3 7 9 10 11 5

Sorted suffixes:

15 \$
14 A\$
0 AACGATAGCGGTAGA\$
1 ACGATAGCGGTAGA\$
12 AGA\$
6 AGCGGTAGA\$
4 ATAGCGGTAGA\$
2 CGATAGCGGTAGA\$
8 CGGTAGA\$
13 GA\$
3 GATAGCGGTAGA\$
7 GCGGTAGA\$
9 GGTAGA\$
10 GTAGA\$
11 TAGA\$
5 TAGCGGTAGA\$

شکل ۸: نمونه سوم

```

1 using System;
2 using TestCommon;
3
4 namespace A6
5 {
6     public class Q4ConstructSuffixArray : Processor
7     {
8         public Q4ConstructSuffixArray(string testDataName)
9         : base(testDataName) { }
10
11         public override string Process(string inStr) =>
12             TestTools.Process(inStr, (Func<String, long[]>)Solve);
13
14         /// <summary>
15         /// Construct the suffix array of a string
16         /// </summary>
17         /// <param name="text"> A string Text ending with a "$" symbol </param>
18         /// <returns> SuffixArray(Text), that is, the list of starting positions
19         /// (0-based) of sorted suffixes separated by spaces </returns>
20         public long[] Solve(string text)
21         {
22             throw new NotImplementedException();
23         }
24     }
25 }

```