

حذف لایه‌های جدید اضافه شده به فایل‌های PDF توسط iTextSharp

عنوان:

وحید نصیری

نویسنده:

۸:۰ ۱۳۹۲/۰۴/۳۱

تاریخ:

www.dotnettips.info

آدرس:

برچسب‌ها: iTextSharp, Watermark

شاید یک سری از Ebook‌های PDF ایی را دیده باشید که سایت‌های ثالث، آن‌ها را پس از افزودن لایه‌ای متنی، مثلاً در ذیل تمام صفحات به همراه آدرس وب سایت خودشان، باز انتشار می‌دهند. در مطلب جاری قصد داریم، نحوه حذف این لایه‌های اضافی را توسط iTextSharp بررسی کنیم.



MANNING

\$49.99 / Can \$52.99 [INCLUDING eBook]

www.-ebooks.info

یافتن و حذف لایه‌های اضافه شده به صفحات یک فایل PDF

برای آشنایی با ساختار سطح پایین لایه‌های اضافه شده نیاز است به برنامه iText Rups مراجعه کنیم.

The screenshot displays the iTextSharp PDF viewer interface. The top pane shows the document structure tree, and the bottom pane shows the stream content of the selected page (Page 1).

Document Structure Tree:

- /Kids: [15626 0 R, 15627 0 R, 15628 0 R, 15629 0 R, 15630 0 R]
- 15626 0 R -> Dictionary of type: /Pages
 - Dictionary of type: /Pages
 - /Parent: 15625 0 R -> Dictionary of type: /Pages
 - /Type: /Pages
 - /Count: 51
 - /Kids: [596 0 R, 620 0 R, 727 0 R, 837 0 R, 887 0 R, 929 0 R]
 - 596 0 R -> Dictionary of type: /Pages
 - Dictionary of type: /Pages
 - /Parent: 15626 0 R -> Dictionary of type: /Pages
 - /Type: /Pages
 - /Count: 6
 - /Kids: [15631 0 R, 8008 0 R, 595 0 R, 599 0 R, 611 0 R, 616 0 R]
 - 15631 0 R -> Dictionary of type: /Page
 - Page 1
 - /CropBox: [0, 0, 531, 666]
 - /Parent: 596 0 R -> Dictionary of type: /Pages
 - /Contents: 16932 0 R, 15839 0 R, 16930 0 R, 16933 0 R
 - 16932 0 R -> Stream
 - Stream
 - 15839 0 R -> Stream
 - Stream
 - 16930 0 R -> Stream
 - Stream
 - 16933 0 R -> Stream
 - Stream
 - 16934 0 R -> Stream
 - Stream

Stream Content of Page 1:

Key	Value
/QQAP	0 Ts
178	0 Tw 0 Tc
	BT
	1 0 0 1 0 0 Tm
	(www. ebooks.info)Tj
	ET
	1 0 0 1 0 0 cm
	1 0 0 1 -223 -7 cm

همانطور که مشاهده می‌کنید، برای رسیدن به لایه‌ای که حاوی متن اضافه شده به ذیل تمام صفحات است، نیاز است ابتدا صفحات را گشوده و سپس CONTENTS آن‌ها را استخراج کنیم. در این CONTENTS کلیه stream‌های موجود را بررسی و هر کدام که حاوی متن مورد نظر ما بودند، یافته و سپس آن استریم را با مقدار دهی طول آن به صفر، حذف کنیم. روش کار را در متد ذیل مشاهده می‌کنید:

```
private static void removeWatermarkLayer(string watermarkedFile, string text, string
unwatermarkedFile)
{
    PdfReader.unethicalreading = true;
    PdfReader reader = new PdfReader(watermarkedFile);
    reader.RemoveUnusedObjects();
    int pageCount = reader.NumberOfPages;
    for (int i = 1; i <= pageCount; i++)
    {
        var page = reader.GetPageN(i);
        var contentarray = page.GetAsArray(PdfName.CONTENTS);
        if (contentarray == null)
            continue;

        for (int j = 0; j < contentarray.Size; j++)
        {
            var stream = (PRStream)contentarray.GetAsStream(j);
            //دریافت محتوای خام صفحه
            var content =
System.Text.Encoding.ASCII.GetString(PdfReader.GetStreamBytes(stream));
            if (content.Contains(text))
            {
                //حذف کامل محتوا از فایل
                stream.Put(PdfName.LENGTH, new PdfNumber(0));
                stream.SetData(new byte[0]);
            }
        }
    }

    using (var fileStream = new FileStream(unwatermarkedFile, FileMode.Create,
FileAccess.Write, FileShare.None))
    {
        using (var stamper = new PdfStamper(reader, fileStream))
        {
            {
                stamper.SetFullCompression();
                stamper.Close();
            }
        }
    }
}
```

در این متد همان watermarkedFile فایل اصلی دارای لایه‌های اضافی است. Text متنی است که در استریم‌های صفحات به دنبال آن خواهیم گشت و unwatermarkedFile نام و مسیر فایل تصحیح شده نهایی است که قرار است تولید شود.