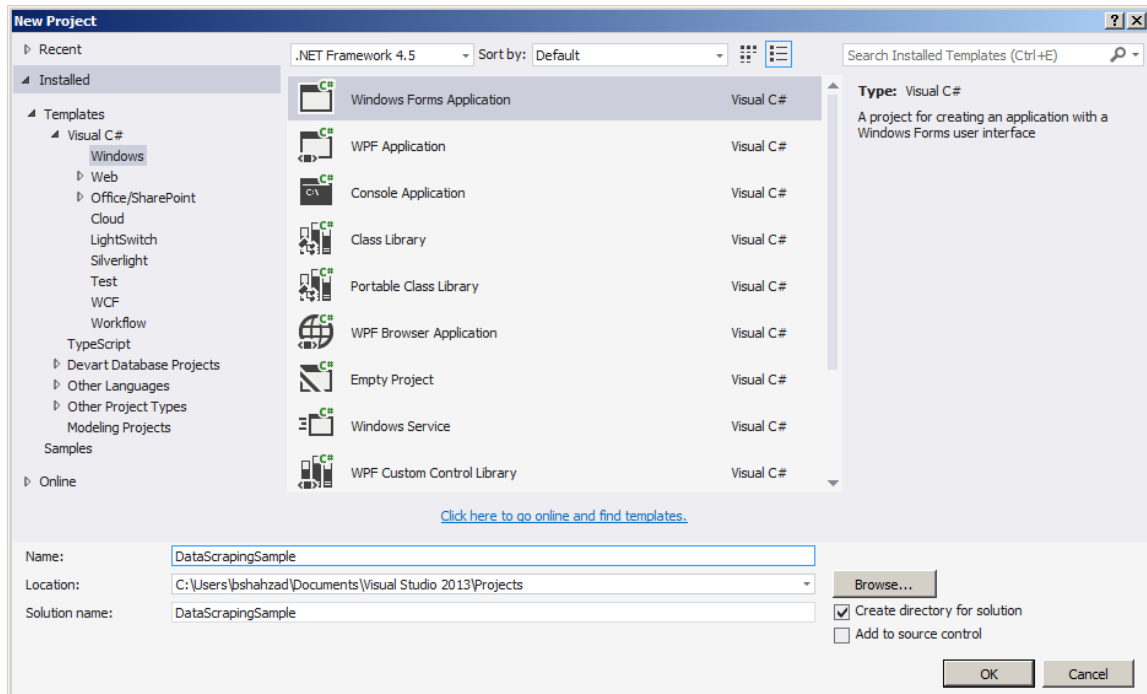# Data scraping using .NET Framework Classes + Fizzler

This guide will explain how to use .NET classes to get HTML response against a request. Once HTML is received then you can use different libraries (e.g. Fizzler, CsQuery). Here we've used Fizzler +HTMLAgilityPack. This tutorial will guide you about adding these packages using NuGet. It explains how to extract information from a web page. This tutorial is only for learning purposes.

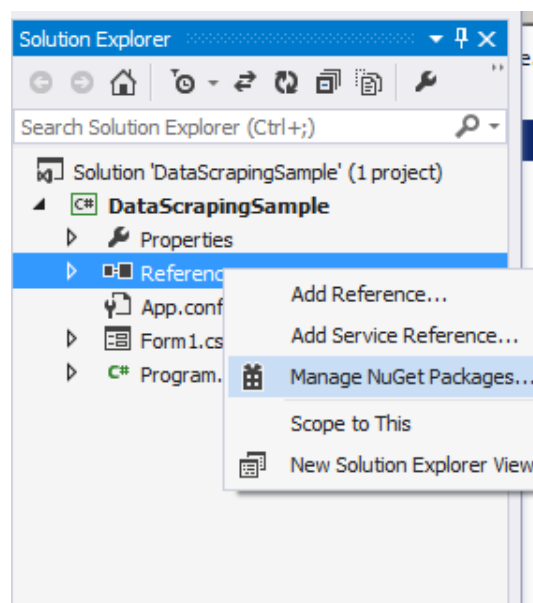| Version | Last updated | Comments | Modified By |
|---------|--------------|----------|-------------|
| V1.0 | 22-04-2016 | | Bilal Shahzad |

Here we are going to build a simple windows application which will extract basic job information from "rozee.pk" site.
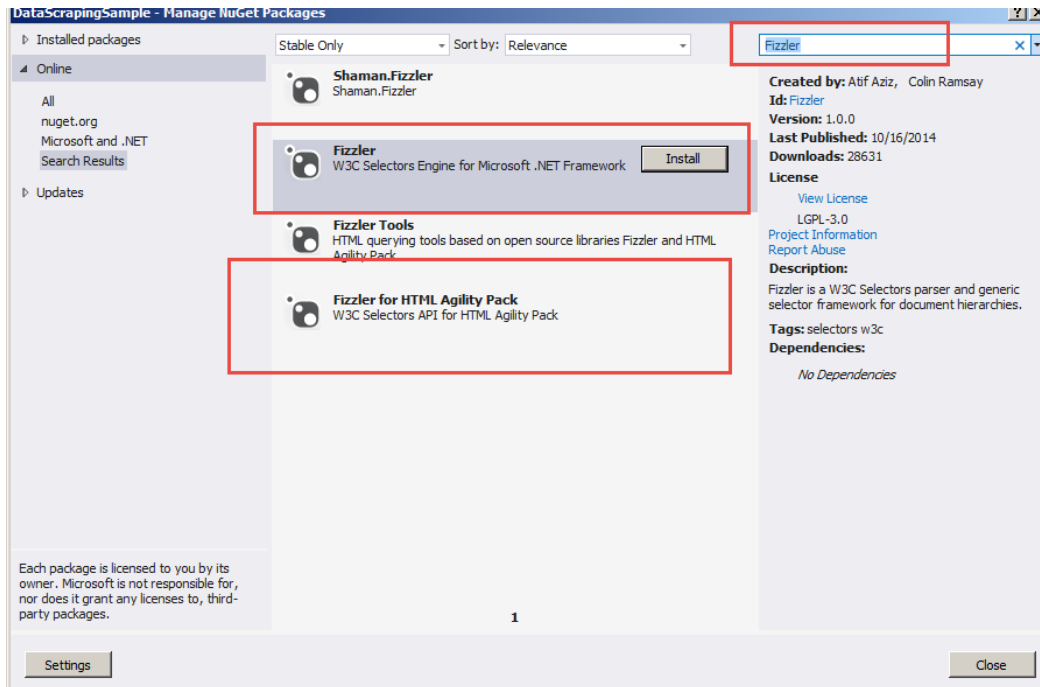
## Walkthrough

1- Create a new "Windows Forms" application with .NET Framework 4.5.
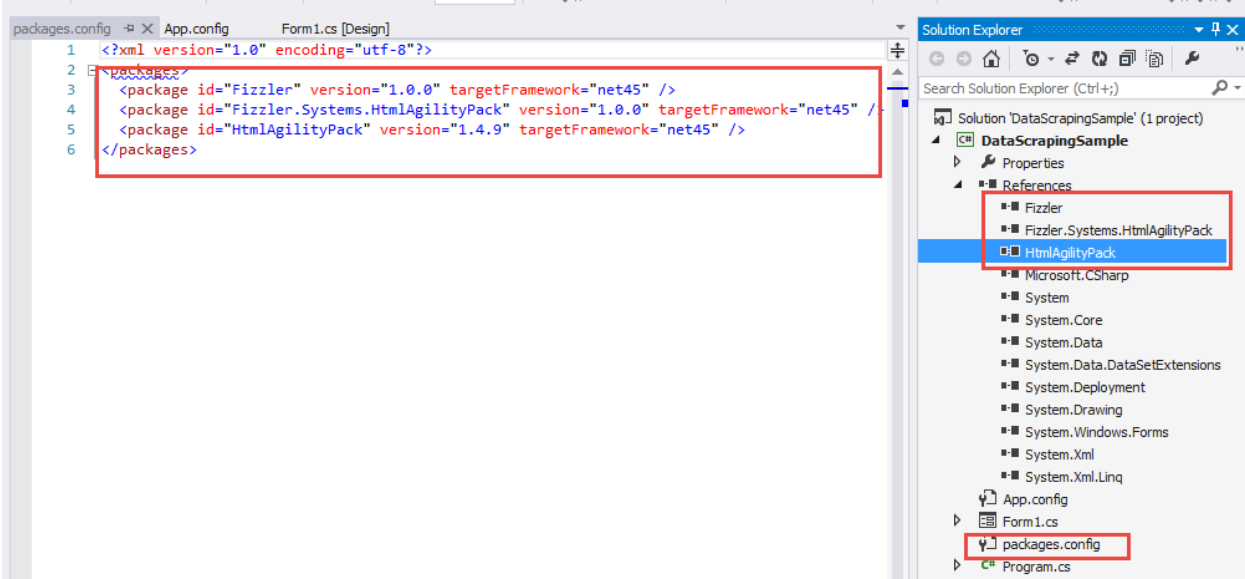


2- Now we need to add some references to access "Fizzler + HtmlAgilityPack" APIs. Right click on "References" and select "Manage NuGet Packages…"
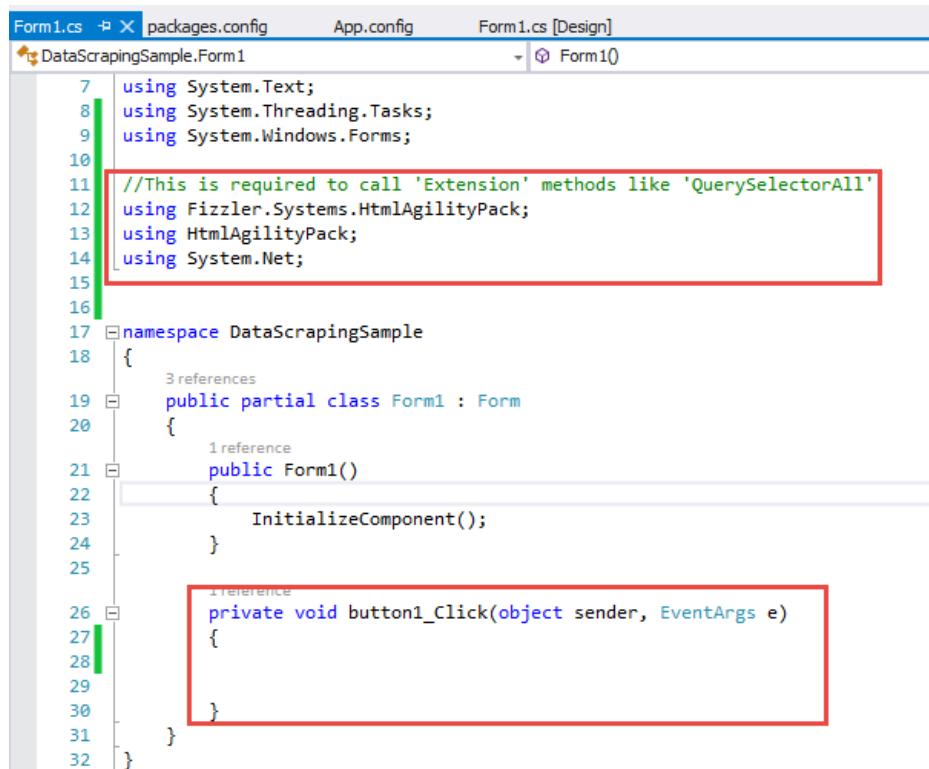
3- In search box type "Fizzler" and Install the packages as highlighted in following screenshot.



4- In "Packages.config", you will see following packages are added. Also check the reference of DLLs shown in right side panel.
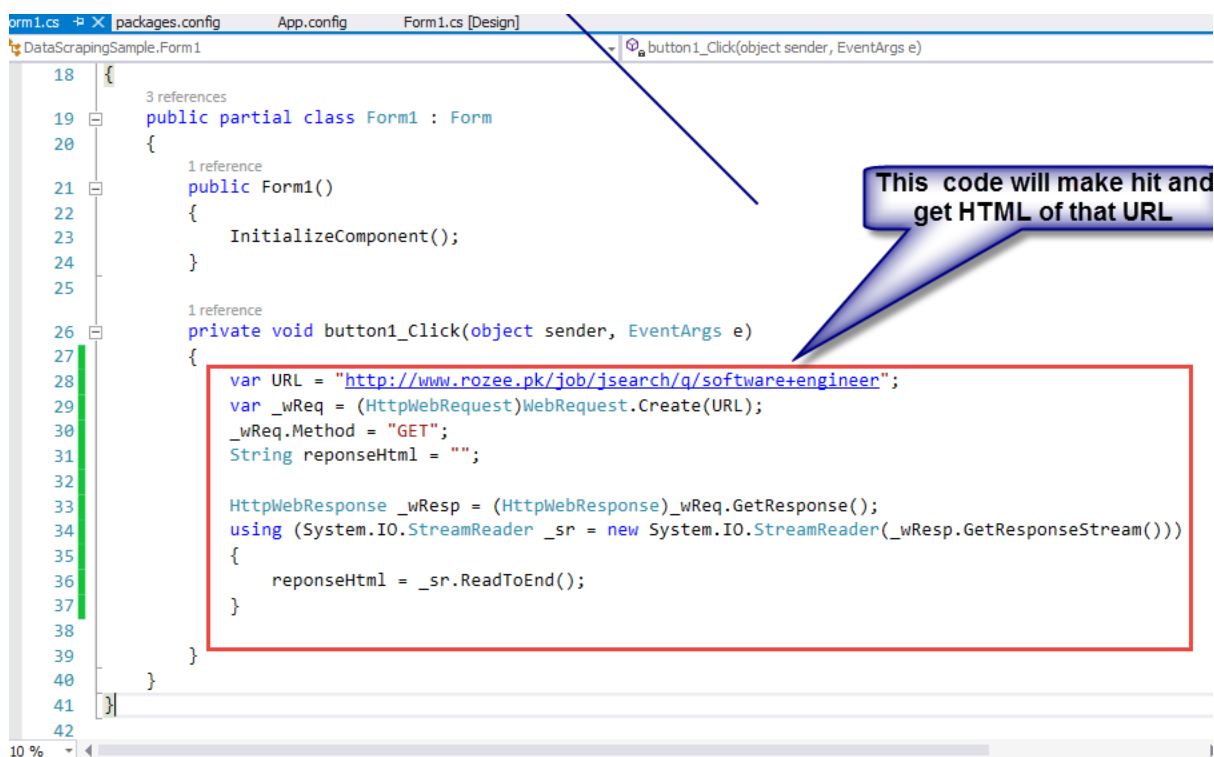
5- Now add a button on a form and generate click event of it. In Code, add the reference of namespaces as shown below.

```
Form1.cs  ⊣ ✕  packages.config      App.config       Form1.cs [Design]
 ⬥ DataScrapingSample.Form1                              ⌄  ⊙ Form1()
      7    using System.Text;
      8    using System.Threading.Tasks;
      9    using System.Windows.Forms;
     10
     11    //This is required to call 'Extension' methods like 'QuerySelectorAll'
     12    using Fizzler.Systems.HtmlAgilityPack;
     13    using HtmlAgilityPack;
     14    using System.Net;
     15
     16
     17  ⊟ namespace DataScrapingSample
     18    {
              3 references
     19  ⊟      public partial class Form1 : Form
     20         {
                  1 reference
     21  ⊟          public Form1()
     22             {
     23                 InitializeComponent();
     24             }
     25
                  1 reference
     26  ⊟          private void button1_Click(object sender, EventArgs e)
     27             {
     28
     29
     30             }
     31         }
     32    }
```

6- We are going to scrap data from following http://www.rozee.pk/job/jsearch/q/all. If you type "Software Engineer" in highlighted textbox below and press highlighted search button, page will be redirected to http://www.rozee.pk/job/jsearch/q/software+engineer. We want to scrap this page.

```
orm1.cs  ⊣ ✕  packages.config      App.config       Form1.cs [Design]
 ⬥ DataScrapingSample.Form1                              ⌄  ⊙ button1_Click(object sender, EventArgs e)
     18   {
              3 references
     19  ⊟      public partial class Form1 : Form
     20         {
                  1 reference
     21  ⊟          public Form1()
     22             {
     23                 InitializeComponent();
     24             }
     25
                  1 reference
     26  ⊟          private void button1_Click(object sender, EventArgs e)
     27             {
     28                 var URL = "http://www.rozee.pk/job/jsearch/q/software+engineer";
     29                 var _wReq = (HttpWebRequest)WebRequest.Create(URL);
     30                 _wReq.Method = "GET";
     31                 String reponseHtml = "";
     32
     33                 HttpWebResponse _wResp = (HttpWebResponse)_wReq.GetResponse();
     34                 using (System.IO.StreamReader _sr = new System.IO.StreamReader(_wResp.GetResponseStream()))
     35                 {
     36                     reponseHtml = _sr.ReadToEnd();
     37                 }
     38
     39             }
     40         }
     41   }
     42
 10 %  ⌄ ◂
```
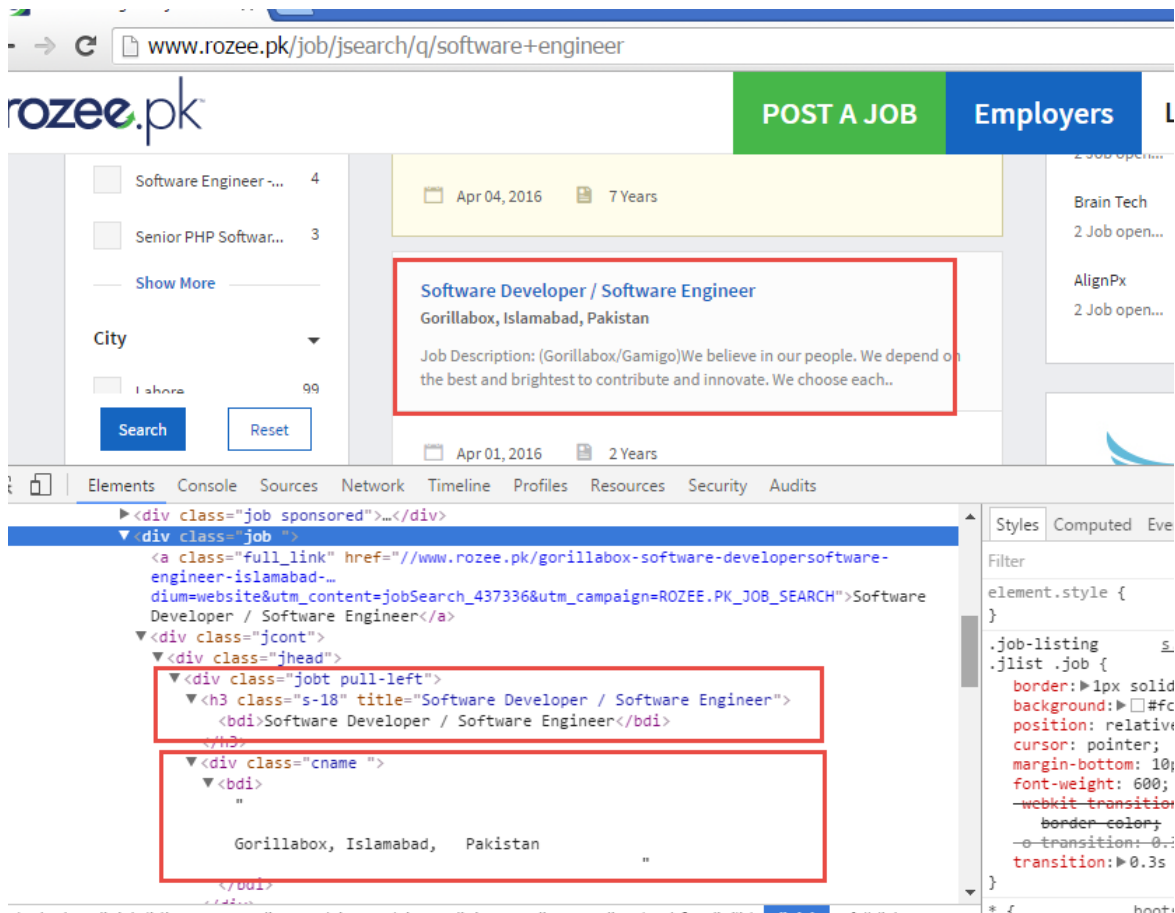
This code will make hit and get HTML of that URL

7- If we open above link in browser, we can investigate the HTML of our results. Each result is in "div.job". Inside div, we can see the selectors of Title and company and so on.



8- Her is the code to find all DIVs and then iterating those divs to extract information from it.



```
28        var URL = "http://www.rozee.pk/job/jsearch/q/software+engineer";
29        var _wReq = (HttpWebRequest)WebRequest.Create(URL);
30        _wReq.Method = "GET";
31        String reponseHtml = "";
32
33        HttpWebResponse _wResp = (HttpWebResponse)_wReq.GetResponse();
34        using (System.IO.StreamReader _sr = new System.IO.StreamReader(_wResp.GetResponseStream()))
35        {
36            reponseHtml = _sr.ReadToEnd();
37        }
38
39        HtmlAgilityPack.HtmlDocument doc = new HtmlAgilityPack.HtmlDocument();
40        doc.LoadHtml(reponseHtml);
41
42        var divs = doc.DocumentNode.QuerySelectorAll(".jlist div.job").ToList();
43
44        foreach (var jobdiv in divs)
45        {
46            var doc1 = new HtmlAgilityPack.HtmlDocument();
47            doc1.LoadHtml(jobdiv.InnerHtml);
48            String jtitle = doc1.DocumentNode.QuerySelector("div.jobt h3 bdi").InnerText;
49
50        }
51
52    }
53 }
```

# Useful Links

http://stackoverflow.com/questions/846994/how-to-use-html-agility-pack

http://www.mikesdotnetting.com/article/273/using-the-htmlagilitypack-to-parse-html-in-asp-net

http://stackoverflow.com/questions/22092208/parsing-html-with-csquery

https://github.com/jamietre/CsQuery