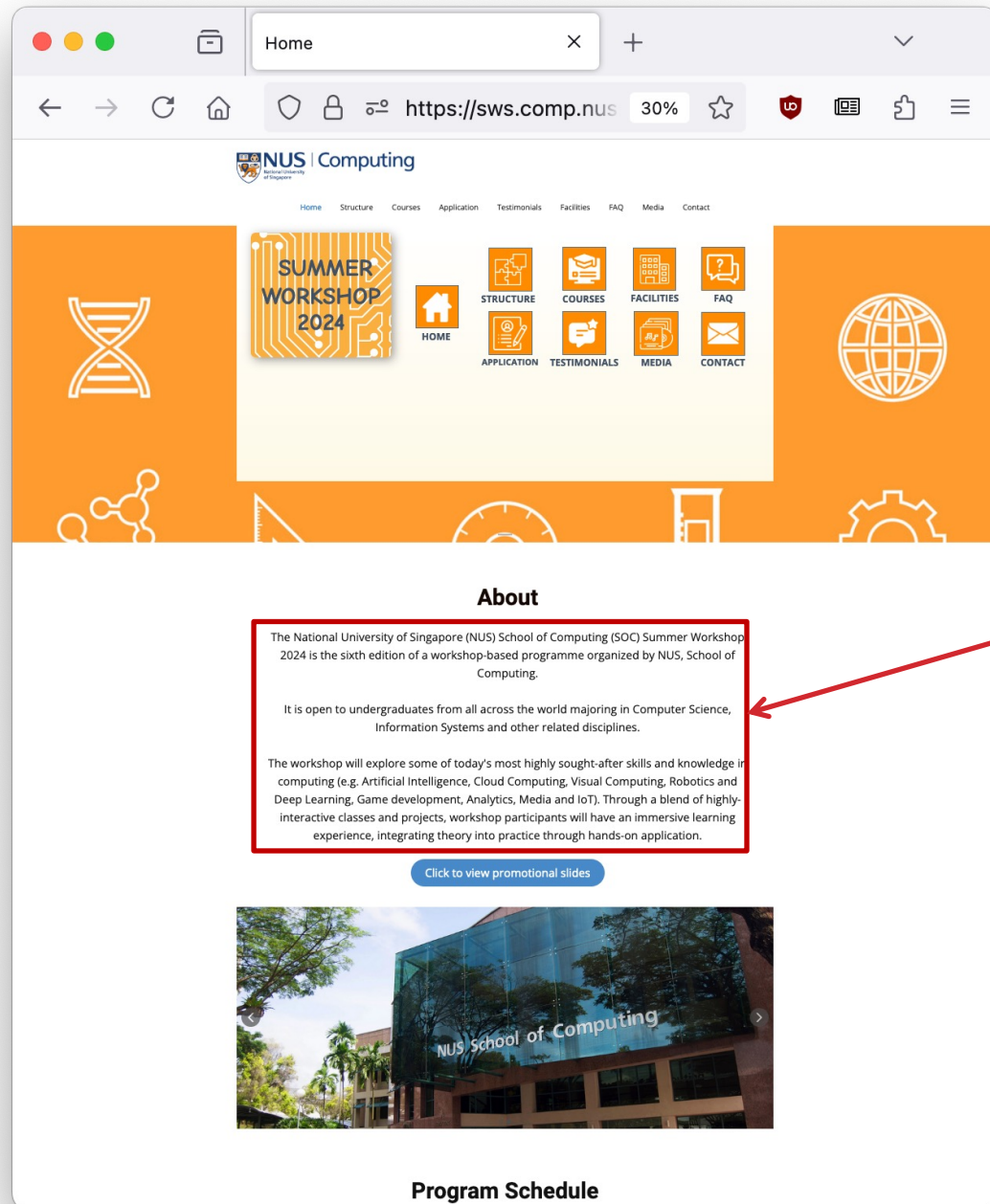# LECTURE 5 MINING WEB CONTENT II

LEK HSIANG HUI

# OUTLINE
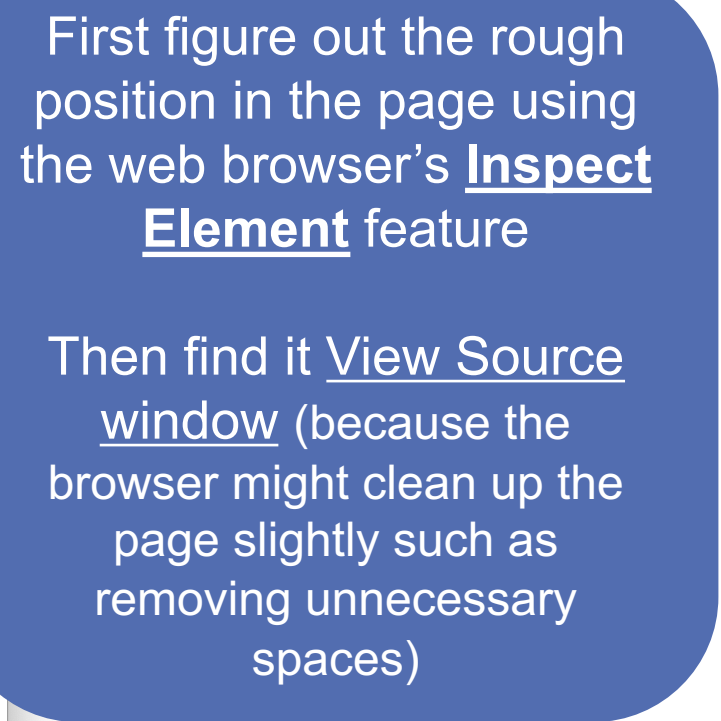
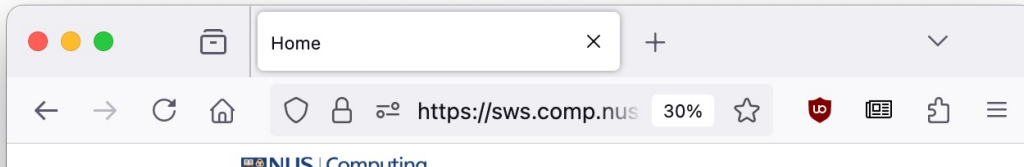**Document Object Model (DOM)**

**XPath**

**CSS Selectors**

**Extracting Content using HTML Parser**

Recap: **Extracting contents from HTML source**

Suppose we want to extract the contents under the About

First figure out the rough position in the page using the web browser's **Inspect Element** feature

Then find it View Source window (because the browser might clean up the page slightly such as removing unnecessary spaces)

Write the regular expression:

`<p class="u-align-center u-text u-text-2">(.*?)</p>`

Problem is that if the web developer were to write the HTML as:

```
<p class="u-align-center  u-text u-text-2">(.*?)</p>
```

If there are additional spaces or different ordering for the CSS classes, the scraper will break

About

computing (e.g. Artificial Intelligence, Cloud Computing, Visual Computing, Robotics and Deep Learning, Game development, Analytics, Media and IoT). Through a blend of highly-interactive classes and projects, workshop participants will have an immersive learning experience, integrating theory into practice through hands-on application.

Click to view promotional slides

NUS School of Computing

Program Schedule

# FLAW OF STRING-BASED APPROACH OF WEB SCRAPING

**Too easily affected by the way how the HTML is written**

- Minor changes (e.g. newlines, spaces, capitalization, shifting of attributes ordering, etc) might break the scraper

- Even when the page is still a totally valid page and might look exactly the same

# RECAP: TECHNIQUES FOR WEB SCRAPING

**The following are some of the techniques for doing web scraping:**

- Extracting content from HTML source
- **Extracting content using a HTML parser**
- Web Scraping using APIs
- Scraping using an actual browser/headless browser

> We will look at another approach which is more robust against this situation

# DOCUMENT OBJECT MODEL (DOM)

| Document Object Model (DOM) | XPath | CSS Selectors | Extracting Content using HTML Parser |

```html
<html>
 <head>
  <meta charset="UTF-8" />
  <meta name="description" content="…"/>
  <meta http-equiv="X-UA-Compatible" content="ie=edge" />
  <title>Welcome to SWS3023</title>
 </head>
 <body>
  <div class="article" id="a0042">
   <h1>Cupcake Article</h1>
   <div class="header">…</div>
   <p>…</p>
   <p>…</p>
  </div>
  <div class="article" id="a0043">
   <h1>Cheese Article</h1>
   <div class="header">…</div>
   <p>…</p>
   <p>…</p>
  </div>
  <div class="article special" id="b0051">
   <h1>Office Article</h1>
   <div class="snippet">…</div>
   <p>…</p>
   <p>…</p>
  </div>
 </body>
</html>
```

Recall: HTML defines the content and the layout of the page



**page.html**

# DOCUMENT OBJECT MODEL (DOM)

**When the browser loads a page, it creates a DOM of the page and use it to render the page**



DOM tree of **page.html**

# VIEWING THE DOM ON WEB BROWSER



Press F12 on Firefox/Chrome and go to the **Inspector** (or **Elements** in Chrome) tab

Or… right click >> Inspect Element

SWS3023 Web Mining

# VIEWING THE DOM ON WEB BROWSER



First Choose the Pick Element Icon

Then hover over the elements on the page
The element will be highlighted at the Inspector tab

13

# VIEWING THE DOM ON WEB BROWSER



Alternatively, could just select an element at the Inspector tab, the selected element will be highlighted on the page

# MODIFYING THE DOM ON WEB BROWSER



Possible to modify the contents of the DOM by double clicking and editing

# MODIFYING THE DOM ON WEB BROWSER



Right click to get more options

Common options:
- Edit HTML
- Delete Node
- Copy (to be elaborated later)

# NAVIGATING THE DOM TREE

**Possible to navigate the DOM programmatically:**

- e.g. <u>extract content</u>, add/remove elements, etc

**2 ways to navigate the DOM tree**

- Using **XPath**
- Using **CSS Selectors**

# XPATH

Document Object Model (DOM) → **XPath** → CSS Selectors → Extracting Content using HTML Parser

# XPATH

**XPath is a way to navigate through elements and attributes in an eXtensible Markup Language (XML) document**

- HTML is very similar to a XML document
- Thus, XPath can be used to navigate through HTML documents also

# XPATH EXPRESSIONS

| XPath Expression | Description |
| --- | --- |
| /html | Select the **html** node |
| /html/head | Select the **head** node (notice the navigation path) |
| /html/body/div | Select all the **div** <u>directly</u> under **body** node |
| /html/body/div[1] | Select the <u>first</u> **div** <u>directly</u> under **body** node |
| //div | Select all **div** in the document (regardless of its ancestor) |
| //div[1] | Select <u>all the first</u> **div** in the document (originating from an ancestor). Note that this can select <u>multiple div</u> |
| //div[last()] | Select <u>all the last</u> **div** in the document (originating from an ancestor). Note that this can select <u>multiple div</u> |
| //body/div | Select all the **div** nodes that is a **direct child** of the **body** node |
| //body//div | Select all the **div** nodes that is a **descendent** of the **body** node |
| //body//* | Select all the nodes under the **body** node (* is a wildcard) |

# TRYING XPATH ON THE BROWSER



1. Select the **Console** tab

3. Result of the node selection is return as an array
(hovering over a value will highlight the element on the page)

2. Type `$x(XPATH_EXPRESSION)`
(this is effectively JavaScript)

# XPATH EXPRESSIONS

**Exercise:** Try out these expression using the browser **Console**

| XPath Expression | Description |
|---|---|
| /html | Select the **html** node |
| /html/head | Select the **head** node (notice the navigation path) |
| /html/body/div | Select all the **div** <u>directly</u> under **body** node |
| /html/body/div[1] | Select the <u>first</u> **div** <u>directly</u> under **body** node |
| //div | Select all **div** in the document (regardless of its ancestor) |
| //div[1] | Select <u>all the first</u> **div** in the document (originating from an ancestor). Note that this can select <u>multiple div</u> |
| //div[last()] | Select <u>all the last</u> **div** in the document (originating from an ancestor). Note that this can select <u>multiple div</u> |
| //body/div | Select all the **div** nodes that is a **direct child** of the **body** node |
| //body//div | Select all the **div** nodes that is a **descendent** of the **body** node |
| //body//* | Select all the nodes under the **body** node (* is a wildcard) |

# XPATH EXPRESSIONS

| XPath Expression | Description |
| --- | --- |
| `//div[@id]` | Select all **div** node with **id** attribute |
| `//div[@class='article']` | Select all **div** node with **class** attribute with the value **article** (must be **exact match**, i.e. even with extra spaces will not work) |
| `//*[@class='article']` | Select **any** node with **class** attribute with the value **article** |
| `//div[starts-with(@id,'a004')]` | Select all **div** node with **id** attribute having value starting with **a004** |
| `//div[contains(@id,'00')]` | Select all **div** node with **id** attribute having value containing **00** |
| `//*[contains(@class,'article')]/..` | Select **parent** node of **any** node with **class** attribute having value containing **article** |
| `//div \| //p` | Select all **div** and **p** nodes |

# CASE SENSITIVITY

| XPath Expression | Case Sensitive? |
|---|---|
| `//div[@id]`<br>`//DIV[@id]`<br>`//div[@ID]`<br>`//DiV[@ID]` | All these are equivalent. Both the node name and the attribute names are **not case sensitive** |
| `//div[@class='article']`<br>`//div[@class='Article']` | These are not equivalent. The attribute value is **case sensitive**<br><br>Same principle applies to `starts-with()` and `contains()` |

# XPATH REFERENCES

[https://www.w3schools.com/xml/xpath_syntax.asp](https://www.w3schools.com/xml/xpath_syntax.asp)

# CSS SELECTORS

Document Object Model (DOM) → XPath → **CSS Selectors** → Extracting Content using HTML Parser

# CSS SELECTORS

**CSS selectors is another way to select elements in HTML**

- Used mainly for selecting HTML elements in order to apply styling into the webpage
- But could also be used for referencing to elements in JavaScript or for doing web scraping
- More commonly used compared to XPath
- Tends to be shorter compared to XPath

# CSS SELECTORS

**CSS selectors is another way to select elements in HTML**

- …
- For the purpose of web scraping, these are the most useful type of selectors
    - **Element** selector
    - **Class** selector
    - **Id** selector
    - **Attribute** selector
    - **Pseudo-Classes** selector
    - **Relationship** selector

# ELEMENT, CLASS, ID SELECTOR

**Element selector used to select all elements with a certain type of tag**

**Id selector used to select the element with a certain id value**

- HTML usually use `id` to uniquely identify an element

**Class selector used to select the element with a certain class attribute value**

- HTML elements usually use `class` to denote the style of the element

# TRYING CSS SELECTORS



To use css selector, type
`$$(CSS_SELECTOR)`

**Element selector**:
`p` will select all paragraph elements

**Id selector**:
`#a0042` will select the element with id = a0042

**Class selector**:
`.article` will select the elements with class <u>containing</u> article class

Browser content:

**Modified Article**

Gummies pie dragée pastry lemon drops. Sweet roll bonbon to sweet roll biscuit.

Topping sesame snaps marzipan. Tootsie roll chocolate bar se snaps bear claw danish. Jelly beans sweet roll jujubes caramel

Halvah chocolate oat cake tiramisu topping apple pie lollipop marshmallow.

**Choose Article**

Console output:

```
$$('p')
Array(6) [ p, p, p, p, p, p ]
$$('#a0042')
Array [ div#a0042.article ]
$$('.article')
Array(3) [ div#a0042.article, div#a0043.article, div#b0051.article.special ]
```

# CSS SELECTORS

Possible to have multiple type of selectors

| CSS Selector | Description |
|---|---|
| `div#a0042` | This selects the **div** node with **id = a0042** |
| `div.article` | This selects all the **div** nodes with **class** <u>containing</u> **article** (<u>need not be an exact match</u>, the element can have multiple class values separated by spaces) |
| `.article#a0042` | This selects all nodes with **class** <u>containing</u> **article** and with **id = a0042** |
| `.article.special` | This selects all nodes with **class** containing both **article** and **special** (can be any ordering but be careful that <u>no spaces between the 2 classes</u>) |
| `div,p`<br>`div , p`<br>`div, p` | This selects all the **div** and **p** nodes |

# ATTRIBUTE SELECTORS

| CSS Selector | Description |
|---|---|
| `meta[name]` | This selects the **meta** node with a **name** attribute |
| `meta[name='description']` | This selects the **meta** node with a **name** attribute with the value **description** (exact match) |
| `meta[content~='page']` | This selects the **meta** node with a **content** attribute containing a **page** as a <u>whole word match</u>. Would not match somepage. |
| `meta[content*='learn']` | This selects the **meta** node with a **content** attribute containing **learn**. Matches as long as there is a **substring** of learn. |
| `[class = 'article']` `*[class = 'article']` | This selects all node with **class** attribute that is **article** (exact match) |
| `[id ^= 'a00']` | This selects all node with **id** attribute that <u>starts with</u> **a00** |
| `[id $= '51']` | This selects all node with **id** attribute that <u>ends with</u> **a00** |

# PSEUDO-CLASSES SELECTORS

| CSS Selector | Description |
|---|---|
| `div:first-child` | This selects all **div** nodes that is the 1<sup>st</sup> child of its parent |
| `p:last-child` | This selects all **p** nodes that is the last child of its parent |
| `p:nth-child(3)` | This selects all **p** nodes that is the 3<sup>rd</sup> node of its parent |

# RELATIONSHIP SELECTORS

| CSS Selector | Description |
|---|---|
| `body div` | This selects all **div** nodes that are <u>descendents</u> of **body** |
| `body > div` | This selects all **div** nodes that are <u>direct child</u> of **body** |
| `h1 + *` | This selects all nodes that is a <u>next sibling</u> of **h1** |
| `.article .snippet` | This will find all nodes having **class** containing **article**, and select all <u>descendent</u> nodes having **class** containing **snippet**. Notice that there is a <u>space separating the 2 classes</u> |

Can mix with the other selectors discussed before

# CASE SENSITIVITY

| CSS Selector | Case Sensitive? |
|---|---|
| `DIV`<br>`Div`<br>`DiV` | These are equivalent. Tag names are **not case sensitive** |
| `.article`<br>`.Article` | These are not equivalent. Classes are **case sensitive** |
| `#a0042`<br>`#A0042` | Likewise, these are not equivalent. Ids are **case sensitive** |
| `meta[name]`<br>`meta[Name]` | These are equivalent. The attribute names are **not case sensitive** |
| `meta[name='description']`<br>`meta[name='Description']` | These are not equivalent. Attribute values are **case sensitive** (as principles as class and id) |

# CSS SELECTOR REFERENCES

**https://www.w3schools.com/cssref/css_selectors.asp**

**XPath vs CSS Selectors:**

**https://johnresig.com/blog/xpath-css-selectors/**
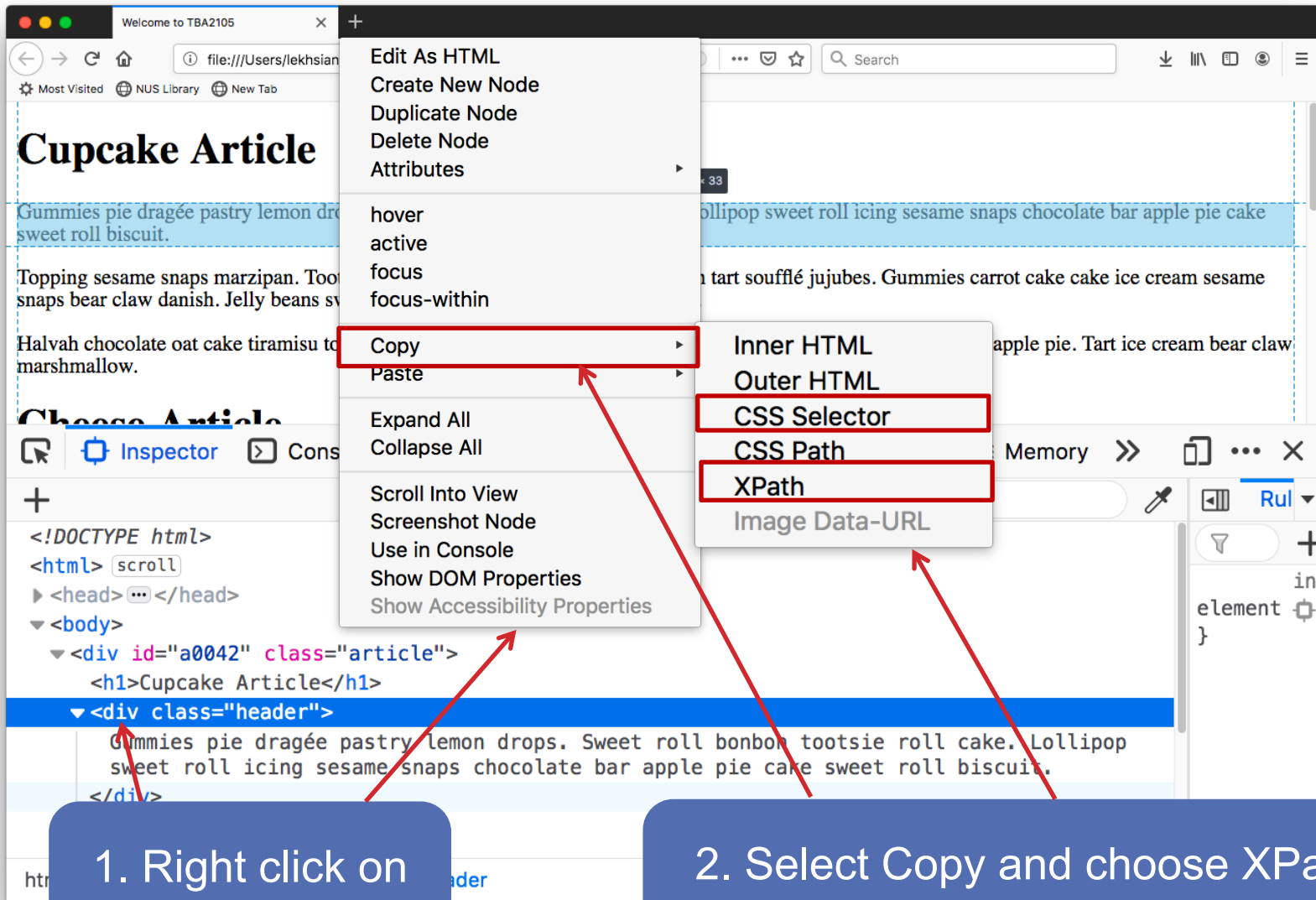
**Not confident of writing your own XPath/CSS Selectors?**

- The browser also allows you to copy the XPath/CSS Selector expressions
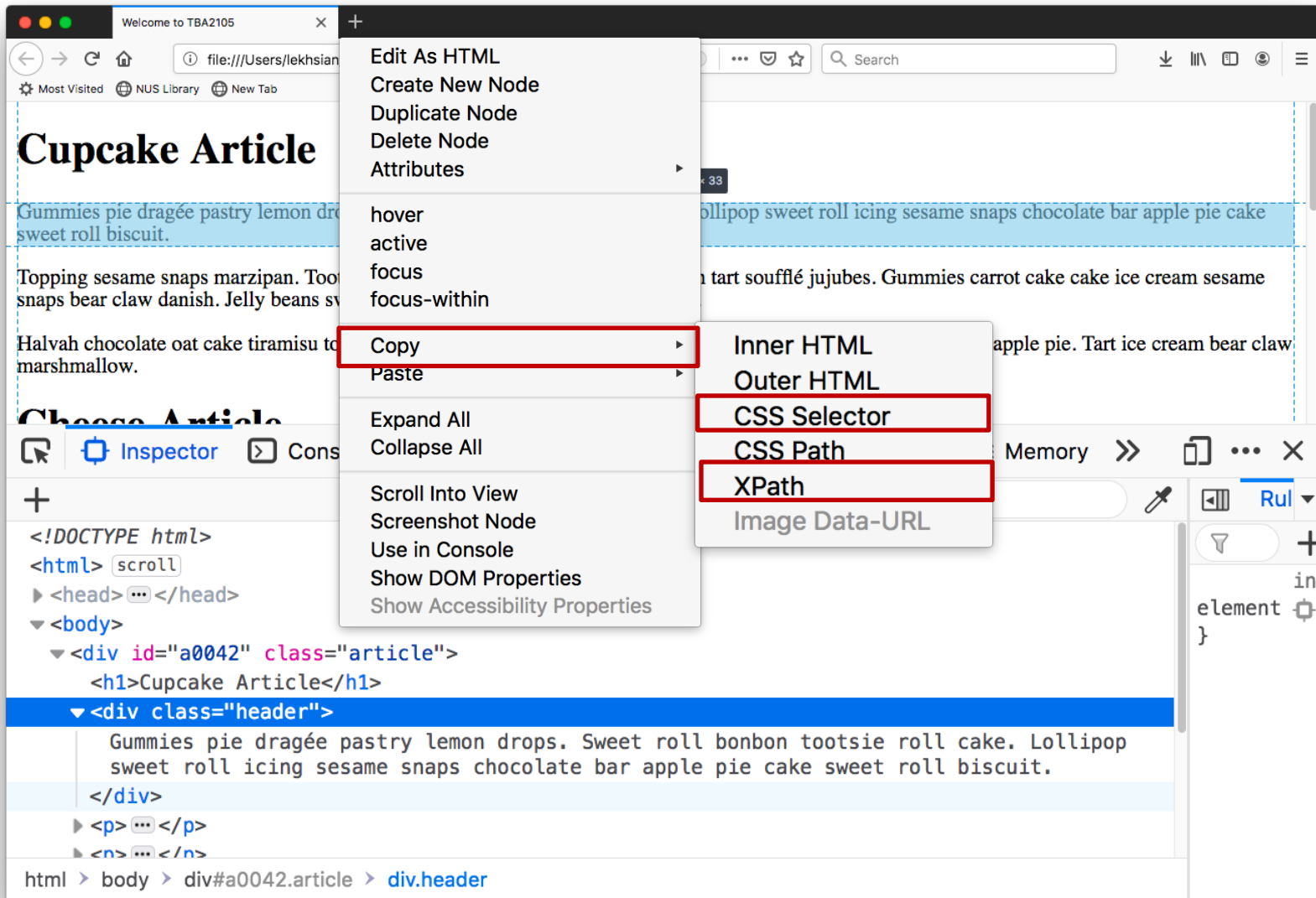
# COPY XPATH/CSS SELECTOR



1. Right click on an element

2. Select Copy and choose XPath or CSS Selector (or Selector for chrome)

# COPY XP[...]

# EXTRACTING CONTENT USING HTML PARSER

| Document Object Model (DOM) | XPath | CSS Selectors | Extracting Content using HTML Parser |
|---|---|---|---|

# HTML PARSING

**Different programming platform has its own library for doing HTML parsing**

- After parsing, we are able to select the HTML elements using XPath and/or CSS selectors
- Can use the **lxml** (for XPath) and **BeautifulSoup** packages

**Idea:**

- Parse the HTML document
- Use XPath/CSS Selector to select the element(s)
- Extract the value (attribute/text/html) of the element(s)

# HANDS-ON: WEB SCRAPING

## XPath approach

```python
from lxml import html
import requests

page = requests.get("https://www.comp.nus.edu.sg/~lekhsian/sws3023/page.html")
page.content
```

```
b'<!DOCTYPE html>\n<html>\n  <head>\n    <meta charset="UTF-8" />\n    <meta\n        name="d
escription"\n        content="This is a dummy page for learning selectors"\n    />\n    <meta
http-equiv="X-UA-Compatible" content="ie=edge" />\n    <title>Welcome to TBA2105</title>\n
</head>\n  <body>\n    <p>this is a paragraph outside any div</p>\n    <div class="article"
id="a0042">\n        <h1>Cupcake Article</h1>\n        <div class="header">\n          Gummies pi
e drag\xc3\xa9e pastry lemon drops. Sweet roll bonbon tootsie roll\n          cake. Lollipop
sweet roll icing sesame snaps chocolate bar apple pie\n          cake sweet roll biscuit.\n
<div>this is a third level div</div>\n        </div>\n        <p>\n          Topping sesame snaps
marzipan. Tootsie roll chocolate bar sesame snaps\n          muffin tart souffl\xc3\xa9 jujub
es. Gummies carrot cake cake ice cream sesame\n          snaps bear claw danish. Jelly beans
sweet roll jujubes caramels cupcake\n          biscuit.\n        </p>\n        <p>\n          Halva
h chocolate oat cake tiramisu topping apple pie lollipop jelly-o\n          cake. Topping cot
ton candy sweet marzipan apple pie. Tart ice cream bear\n          claw marshmallow.\n
</p>\n    </div>\n    <div class="article" id="a0043">\n        <h1>Cheese Article</h1>\n
<div class="header">\n          The big cheese red leicester rubber cheese. Stilton taleggio
halloumi\n          croque monsieur bocconcini cheese triangles cheesecake boursin. Ricotta\n
```

41

# SUMMARY

**Document Object Model (DOM)**

**Navigating the DOM tree**

- XPath and CSS Selectors

**Extracting Content using HTML Parser**

- lxml & BeautifulSoup packages

# WHAT'S NEXT?

**Mining Web Content III**