# Basics of XML

```xml
<?xml version="1.0" encoding="UTF-8"?>
<movies>
    <movie mins="126" lang="eng">
        <title>Good Will Hunting</title>
        <director>
            <first_name>Gus</first_name>
            <last_name>Van Sant</last_name>
        </director>
        <year>1998</year>
        <genre>drama</genre>
    </movie>
    <movie mins="106" lang="spa">
        <title>Y tu mama tambien</title>
        <director>
            <first_name>Alfonso</first_name>
            <last_name>Cuaron</last_name>
        </director>
        <year>2001</year>
        <genre>drama</genre>
    </movie>
</movies>
```

# Basics of XML and HTML

# Goal

## XML & HTML

The goal of these slides is to give you a **crash introduction to XML and HTML** so you can get a good grasp of those formats for the rest of the lectures

# Synopsis

### In a nutshell

We'll cover a the following concepts:

- Importance of XML and HTML
- Hierarchical Structure

# XML and HTML

## Why you should care about XML and HTML?

- Large amounts of data and information are stored, shared and distributed using HTML and XML-dialects
- They are widely adopted and used in many applications
- Working with data from the Web means dealing with HTML

# XML

eXtensible Markup Language

```xml
1 <?xml version="1.0" encoding="ISO8859-1" ?>
2 <CATALOG>
3   <PLANT>
4     <COMMON>Bloodroot</COMMON>
5     <BOTANICAL>Sanguinaria canadensis</BOTANICAL>
6     <ZONE>4</ZONE>
7     <LIGHT>Mostly Shady</LIGHT>
8     <PRICE>$2.44</PRICE>
9     <AVAILABILITY>031599</AVAILABILITY>
10  </PLANT>
11
12  <PLANT>
13    <COMMON>Columbine</COMMON>
14    <BOTANICAL>Aquilegia canadensis</BOTANICAL>
15    <ZONE>3</ZONE>
16    <LIGHT>Mostly Shady</LIGHT>
17    <PRICE>$9.37</PRICE>
18    <AVAILABILITY>030699</AVAILABILITY>
19  </PLANT>
20
21  <PLANT>
22    <COMMON>Marsh Marigold</COMMON>
23    <BOT ANICAL>Caltha palustris</BOTANICAL>
24    <ZONE>4</ZONE>
25    <LIGHT> Mostly   Sunny </LIGHT>
26    <PRICE> $6.81 </PRICE>
27    <AVAILABILITY>   051799 </AVAILABILITY>
28  </PLANT>
29
30  <PLANT>
```

# Some Definitions

*"XML is a markup language that defines a set of rules for encoding documents in a format that is both human-readable and machine-readable"*

http://en.wikipedia.org/wiki/XML

*"XML is a data description language used for describing data"*

Paul Murrell
Introduction to Data Technologies

# Some Definitions

*"XML is a very general structure with which we can define any number of new formats to represent arbitrary data"*

*"XML is a standard for the semantic, hierarchical representation of data"*

Deb Nolan & Duncan Temple Lang

XML and Web Technologies for Data Sciences with R

## XML

XML stands for **eXtensible Markup Language**

## Broadly speaking ...

XML provides a flexible framework to create formats for describing and representing data

# Markups

## Markup

A **markup** is a sequence of characters or other symbols inserted at certain places in a document to indicate either:

- ► how the content should be displayed when printed or in screen
- ► describe the document's structure

## Markup Language

A markup language is a system for **annotating** (i.e. *marking* ) a document in a way that the content is distinguished from its representation (eg LaTeX, PostScript, HTML, SVG)

## XML Markups

In XML (as well as in HTML) the marks (aka *tags*) are defined using angle brackets: **< >**

&lt;mark&gt;Text marked with special tag&lt;/mark&gt;

## Extensible?

The concept of *extensibility* means that we can define our own marks, the order in which they occur, and how they should be processed. For example:

- ► <my_mark>
- ► <awesome>
- ►<boring>
- ►<pathetic>

# About XML

## XML is NOT
- a programming language
- a network transfer protocol
- a database

## XML is
- more than a markup language
- a generic language that provides structure and syntax for representing any type of information
- a meta-language: it allows us to create or define other languages

# Minimalist Example

# XML Example

**Ultra Simple XML**

```
<movie>
  Good Will Hunting
</movie>
```

- ► one single element *movie*
- ► start-tag: `<movie>`
- ► end-tag: `</movie>`
- ► content: Good Will Hunting

# XML Example

## Ultra Simple XML

```
<movie mins="126" lang="en">
   GoodWill Hunting
</movie>
```

- xml elements can have **attributes**
- attributes: mins (minutes) and lang (language)
- attributes are *attached* to the element's start tag
- attribute values **must be quoted!**

# XML Example

**Minimalist XML**

```
<movie mins="126" lang="en">
  <title>Good Will Hunting</title>
  <director>Gus Van Sant</director>
  <year>1998</year>
  <genre>drama</genre>
</movie>
```

- ► an xml element may contain other elements
- ► *movie* contains several elements: *title, director, year, genre*

## Simple XML

```
<movie mins="126" lang="en">
  <title>Good Will Hunting</title>
  <director>
    <first_name>Gus</first_name>
    <last_name>Van Sant</last_name>
  </director>
  <year>1998</year>
  <genre>drama</genre>
</movie>
```

► Now *director* has two child elements: *first name* and *last name*

## Conceptual XML

```
<Root>
  <child_1>...</child_1>
  <child_2>...</child_2>
    <subchild>...</subchild>
  <child_3>...</child_3>
</Root>
```

- ► An XML document can be represented with a **tree structure**
- ► An XML document must have **one single Root** element
- ► The Root may contain child elements
- ► A child element may contain subchild elements