



COMP20008 Elements of Data Processing

Semester 2 2018

Lecture 2: Data formats: structured, unstructured and semi-structured



THE UNIVERSITY OF
MELBOURNE

Announcements

- Student representatives
 - Shui Jen Wong - shuiw1@student.unimelb.edu.au
 - Yao Peter Shi - yaos5@student.unimelb.edu.au

Please contact them with any feedback you have about the subject



THE UNIVERSITY OF
MELBOURNE

Company scenario – TelOptaphone

Your manager calls you into her office. "Next week we are making a presentation to the CEO and we need to include a profile of all our customers: who they are, their past purchasing behaviour and all the types of interactions they've had with our company. We're going to need it by Monday



The data exists, kind of, but it's spread across multiple systems and is in many types of formats: CSV, XML, JSON, HTML, spreadsheets,

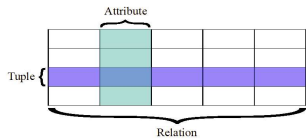


THE UNIVERSITY OF
MELBOURNE

Today's lecture: Data formats

- Where is the data?
- How is data stored and in what formats?
 - *Structured*: Relational databases, CSV
 - *Unstructured*: text
 - *Semi-structured*: HTML, XML, JSON
- Question: Why do we have different data formats and why do we wish to transform between different formats?
- Our purpose for next 2 lectures
 - *To understand differences between and motivation/purposes of these formats*

- It is good to have structure for data!
 - Easier to analyse, easier to query
 - Easier to store
 - Easier to clean, maintain consistency and security, especially with multiple users



- Relational databases, the classic method of storing structured data (banking, sales, airlines ...)
 - Data stored in tables, each row is a data item and columns describe attributes of the data item
 - Can query the data using a high level language such as SQL

Attributes

customer_id	customer_name	customer_street	customer_city	account_number
192-83-7465	Johnson	12 Alma St.	Palo Alto	A-101
192-83-7465	Joluson	12 Alma St.	Palo Alto	A-201
677-89-9011	Hayes	3 Main St.	Harrison	A-102
182-73-6091	Turner	123 Putnam St.	Stamford	A-305
321-12-3123	Jones	100 Main St.	Harrison	A-217
336-66-9999	Lindsay	175 Park Ave.	Pittsfield	A-222
019-28-3746	Smith	72 North St.	Rye	A-201

customer_id	customer_name	customer_street	customer_city
192-83-7465	Johnson	12 Alma St.	Palo Alto
677-89-9011	Hayes	3 Main St.	Harrison
182-73-6091	Turner	123 Putnam Ave.	Stamford
321-12-3123	Jones	100 Main St.	Harrison
336-66-9999	Lindsay	175 Park Ave.	Pittsfield
019-28-3746	Smith	72 North St.	Rye

(a) The customer table

account_number	balance
A-101	500
A-215	700
A-102	400
A-305	350
A-201	900
A-217	750
A-222	700

(b) The account table

customer_id	account_number
192-83-7465	A-101
192-83-7465	A-201
019-28-3746	A-215
677-89-9011	A-102
182-73-6091	A-305
321-12-3123	A-217
336-66-9999	A-222
019-28-3746	A-201

(c) The depositor table

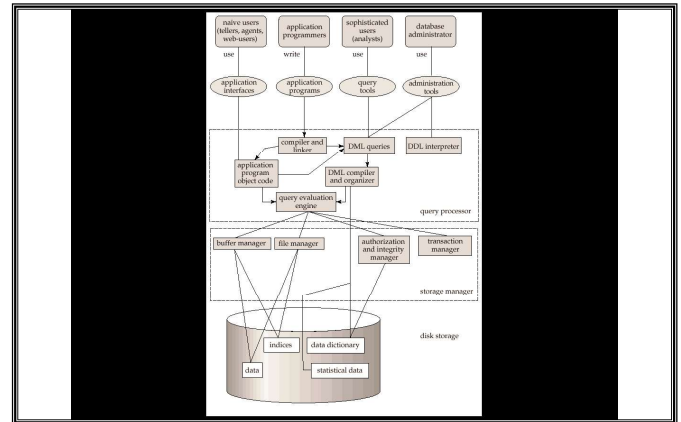
- create table **branch**
 - (branch_name char(15) not null,
 - branch_city char(30),
 - assets integer,
 - primary key (branch_name))



```

select      account.balance
from        depositor, account
where       depositor.customer_id = '192-83-7465'
and         depositor.account_number=account.account_number

```



- In INFO20003 subject you would cover topics like
 - SQL
 - Specification of integrity constraints
 - Data modelling and relational database management systems
 - Transactions and concurrency control
 - Storage management
 - Web-based databases
 -
- Highly relevant to data wrangling!
 - Useful to do INFO20003 as part of a data science specialisation



- Once data is into a relational database, it is easier to wrangle.
 - But maybe hard to load it there in the first place ...
 - Unstructured data: text, HTML - lack regularity

- Huge amounts of data lives in spreadsheets
 - Businesses
 - Hospitals
 -
- Microsoft (Excel), OpenOffice (Calc), Google docs
- CSV (comma separated values) also very popular
 - These are human readable, versus binary XLS format (Excel)
 - CSVs lack the formatting information of an XLS file
- Python libraries
 - csv
 - xlrd, openpyxl
 - **pandas** read_csv function

- Spreadsheets
- Easy to use
- Structured, but not like a relational DB

	A1	B	C	D	E
1	Date	17/02/2014	22/02/2014	28/02/2014	
2	BP	130/80	140/90	131/87	
3	Glucose	4.5	6	5.3	
4	Heart rate	55	70	57	


```

text_example.csv - Notepad
File Edit Format View Help
Date,17/02/2014,22/02/2014,28/02/2014
BP,130/80,140/90,131/87
Glucose,4.5,6,5.3
Heart rate,55,70,57
  
```

Information in tabular format (transactional, simple but many entries)

	A	B	C	D	E	F	G	H	I	J	K	L
1		0-4	5-9	10-14	15-19	20-24	25-29	30-34	35-39	40-44	45-49	50-54
2	1901	434741	456981	434152	378115	350993	320544	292071	272710	221190	155009	119072
3	1911	525633	453246	428161	448536	446270	386376	330960	291432	269518	241616	192919
4	1921	603600	597300	530000	470000	450100	462000	448200	390200	332500	263600	255100
5	1922	611900	602500	546100	482200	456000	457900	459700	402700	344400	287800	261700
6	1923	623200	601800	560600	497300	461500	456800	465200	419800	355600	296300	266800
7	1924	633100	593800	579300	512500	468200	455600	469700	434800	368400	305900	271900
8	1925	642600	590100	596400	529800	480100	465100	472600	446100	379900	317800	274000
9	1926	640800	602200	606100	545700	493400	472500	475000	455900	391600	329000	276100
10	1927	657000	613200	613000	563800	510900	489500	474400	460100	405300	341000	279800
11	1928	634700	626700	613200	580000	530000	496600	475500	474300	421400	351200	286400
12	1929	631600	637500	605100	598000	545300	503000	474300	476600	433600	361900	293200
13	1930	6214	A	B	C	D	E	F	G			
14	1931	6115	1	Rank	Album	Artist	Year	Total Sales	Origin	Genre		
15	1932	6940	2	1	GREATEST HITS	QUEEN	1981	677810	UK	Rock		
16	1933	5740	3	2	SGT. PEPPER'S LONELY HEARTS CLUB BAND	BEATLES	1967	490288	UK	Rock/Pop		
17	1934	5551	4	3	GOLD - GREATEST HITS	ABBA	1992	4610613	Sweden	Pop		
18	1935	5397	5	4	WHAT'S THE STORY MORNING GLORY	OASIS	1996	4415266	UK	Rock		
19	1936	5297	6	5	BROTHERS IN ARMS	DIRE STRAITS	1985	4089764	USA	Rock		
20	1937	5360	7	6	THE DARK SIDE OF THE MOON	PINK FLOYD	1973	3966177	UK	Rock		
21	1938	5452	8	7	THRILLER	MICHAEL JACKSON	1982	3625667	USA	Pop		
22	1939	5590	9	8	GREATEST HITS II	QUEEN	2000	3746404	UK	Rock		
23	1940	5724	10	9	BAD	MICHAEL JACKSON	1987	3545301	USA	Pop		
24	1941	5896	11	10	THE IMMACULATE COLLECTION	MADONNA	1990	3402160	USA	Pop		
25	1942	6102	12	11	STARS	SIMPLY RED	1991	3401082	UK	Pop		
26			13	12	COME ON OVER	SHANIA TWAIN	1998	3368841	Canada	Country/Pop		
27			14	13	RUMOURS	FLEETWOOD MAC	1977	325818	UK	Rock		
28			15	14	BACK TO BLACK	JAMIE BLUNT	2004	317269	UK	Pop		
29			16	15	URBAN HIMALAYS	VERVE	1997	3167976	UK	Pop		
30			17	16	NO ANGEL	DIDO	2003	3045208	UK	Pop		
31			18	17	BRIDGE OVER TROUBLED WATER	SMOKE & GARFUNKEL	1970	3042342	USA	Folk		
32			19	18	BACK TO BLACK	AMY WINEHOUSE	2006	2995303	UK	Retro Soul		
33			20	19	TALK ON CORNERS	THE CORRS	1997	2647666	Ireland	Rock		
34			21	20	BAT OUT OF HELL	MEAT LOAF	1978	2642717	USA	Rock		
35			22	21	SPICE	SPICE GIRLS	1996	2638739	UK	Pop		
36			23	22	WHITE LADDER	DAVID GRAY	2000	2608795	UK	Alternative Rock/Folk		
37			24	23	DIRTY DANCING	ORIGINAL SOUNDTRACK	1987	2602247	UK	Soundtrack		

- Text files...
- No structure
- Harder to index
- Harder to organise
- Lacks regularity and decomposable internal structure
- How can we process/search for information?

```

Untitled - Notepad
File Edit Format View Help
Date = 17/02/14
BP measurement = 130/80
Glucose measurement = 4.5
Heart rate measurement = 55
Date = 22/02/14
BP measurement = 140/90
Glucose measurement = 6.0
Heart rate measurement = 70
Date = 28/02/14
BP measurement = 131/87
Glucose measurement = 5.3
Heart rate measurement = 57
  
```

- Scenario: we have a large collection of unformatted text data. You need to write wrangling code in order to
 - Check if it contains any IP addresses (e.g. 128.250.65.5)
 - Find all the IP addresses
- Requirements
 - Do it succinctly
 - Do it unambiguously
 - Have maintainable code

- Specifying patterns in text – **regular expressions**
 - Good for computing statistics, checking integrity, filtering, substitutions
- Specifying patterns in text
 - '.' matches any character
 - '^' matches start of string
 - '\$' matches end of string
 - '*' zero or more repetitions
 - '+' one or more repetitions
 - '|' the "or" operator, used in conjunction with parentheses ()
 - '[' a set of characters, e.g. [abcd] or [a-zA-Z]
- <https://docs.python.org/2/howto/regex.html>
- regex101.com

- '.' matches any character
- '^' matches start of string
- '\$' matches end of string
- '*' zero or more repetitions
- '+' one or more repetitions
- '|' the "or" operator, used in conjunction with parentheses ()
- '[' a set of characters, e.g. [abcd] or [a-zA-Z]
- JA The pattern "J" followed immediately by 'A'
- JA* The pattern: 'J' followed by **zero** or more occurrences of 'A'
- (J|A)* Zero or more repetitions of 'J' or 'A's
- (J|A)+ One or more repetitions of 'J' or 'A's

- Write regular expressions to specify each of the following
 - Two occurrences of letter 'e' followed immediately by one 'n' and then at least one 't'
 - An 'h' or an 'e' or an 'x', followed by at least one 'a', followed by an 'r'
 - Any 3 characters, possibly followed by a repeated sequence of the character 'x', followed by a 'c' or a 'd'



- What do you think this pattern is for?
 - `[a-zA-Z0-9_+~]@[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+`
 - Could it be improved?



Another type of data

HTML



- Marked up with *elements*, delineated by start and end *tags*. Elements correspond to logical units, such as a heading, paragraph or itemised list.
- **Tags:** Keywords contained in pairs of angle brackets.
 - Not case sensitive.
- Browser determines how to display/present the logical units
- Not all elements need both start and end tags.
- Some elements can have *attributes*. Ordering of attributes is not significant.



```
<div class="icon section5">
<h2><a href="about/index.html">About the Melbourne School of Engineering</a></h2>
<ul>
<li><a href="about/dean_welcome.html">Dean's Welcome</a></li>
<li><a href="about/staff.html">Leadership & Professional Staff</a></li>
<li><a href="about/contact.html">Contact Us</a></li>
<li><a href="http://www.ecr.unimelb.edu.au">ECR: Computer Resources</a></li>
<li><a href="intranet/index.html">For Staff (intranet)</a></li>
<li><a href="casual_staff/index.html">For Casual Staff</a></li>
<li><a href="intranet/review/prof_staff.html">Professional Staff Review</a></li>
<li><a href="/about/safety/index.html">Environment, Health & Safety</a></li>
<li><a href="/about/committees/index.html">Committees</a></li>
</ul>
```

Try it yourself: https://www.w3schools.com/html/tryit.asp?filename=tryhtml5_browsers_myhero

HTML examples: https://www.w3schools.com/html/html_lists.asp

- HTML was designed for pure presentation
- **HTML is concerned with formatting not meaning**
 - it doesn't matter what it is about, HTML will format it
- HTML is not extensible
 - can't be modified to meet specific domain knowledge
 - browsers have developed their own tags (*<bgsound>*, *<layer>*)
- HTML can be inconsistently applied
 - almost everything is rendered somehow
 - e.g. *is this acceptable?</i>*

- Developed in the mid 90's by committee
- Derived from SGML
- A 'meta' mark-up language
 - Used to create other mark up languages
 - Extensible, *user defined tags*
- Separates style and content
- Supports internationalisation (Unicode)
- Rigorous adherence to rules
- Device and system *independent*
- Applications may generate and process XML
- Enables data exchange between different platforms
- Facilitates better encoding of *semantics*
- **Both humans and machines can read it...**
- "Transcends politics through sheer usefulness..."
 - "Intro to XML", Tim Anderson, 2004
 - <http://www.itwritng.com/xmlintro.php>

Hamlet: Act one

SCENE ONE: *Elsinore. A terrace in front of a castle. Francisco is on sentinel duty. Enter Bernardo*

BERNARDO: Who's there?

FRANCISCO: Nay, answer me. Stand and unfold yourself

```
<body>
  <h1> Act One </h1>

  <p>
    SCENE ONE: <i> Elsinore. A terrace in front of a castle.
    Francisco is on sentinel duty. Enter Bernardo. </i>
  </p>
  <p>
    <b> BERNARDO: </b> Who's there?
  </p>
  <p>
    <b> FRANCISCO: </b> Nay, answer me: stand, and
    unfold yourself
  </p>
</body>
```

```
<?xml version="1.0"?>
<act>
  <title> Act One </title>
  <scene>
    <title> SCENE ONE </title>
    <location> Elsinore. A terrace in front of a castle. </location>
    <stagedir> Francisco is on sentinel duty. Enter Bernardo </stagedir>

    <speech>
      <speaker> BERNARDO </speaker>
      <line> Who's there? </line>
    </speech>
    <speech>
      <speaker> FRANCISCO </speaker>
      <line> Nay, answer me: stand, and unfold yourself. </line>
    </speech>
  </scene>
</act>
```

- Can be viewed as a tree of nodes: parent and leaf nodes
 - Try it yourself: <https://codebeautify.org/xmlviewer#>

- xml files must begin with declaration
 - `<?xml version="1.0"?>`
- xml files must have one single root element
 - E.g. `<act>...</act>`
- elements are built with tags, must be properly closed
 - opening `<firstname>` and closing `</firstname>`
 - empty `
`
- an element may have one or more attributes, attributes must be in quotes
 - `<person title="Sir">Richard</person>`
 - `<person title="Mr" sex="Male">James</person>`

- xml code is case sensitive
 - `<title>` is not the same as `<Title>`
- elements must be appropriately nested
 - `<author><firstname>James</firstname></author>`
 - `<author><firstname>James</author> </firstname>`
 - Wrong...
- comments
 - `<!-- comments do not affect the document,`
 - `it's not part of the data that you want to represent -->`

- some characters have special meaning
 - `<` and `&` are strictly illegal inside an element
 - `<text>all books & videos are now < AUD 10</text>`
 - Wrong...
 - `<text>all books & videos are now < AUD 10</text>`
- CDATA (character data) section may be used inside XML element to include large blocks of text, which may contain these special characters such as `&`, `>`
 - `<![CDATA [... ...]]>`
 - `<![CDATA [all books & videos are now < AUD 10]]>`



```
<?xml version="1.0"?>
```

```
<catalog>
```

```

- <book isbn="1-23456-789-0">
  • <title>Beyond the Clouds</title>
  • <author>
    - <firstname>Rebecca</firstname>
    - <surname>Skye</surname>
    - <picture source="rebecca.jpg" />
  • </author>
- </book>
- <book isbn="0-98765-432-1">
- <title>The Final Straw</title>
- <author>
  • <firstname>James</firstname>
  • <surname>Last</surname>
  • <picture source="james.jpg" />
- </author>
- </book>

```

```
</catalog>
```

- Declaration
- One root element
- Attributes in quotes
- Empty tag
- Opening/closing tags
- Tags correctly nested

- **"WELL FORMED"**



- Given the following data: Yellow Balloon, \$99.99
 - i) What are three possible XML encodings of the balloon ?
 - ii) What are some of the circumstances in which one encoding might be better than the others ?



- Here is some information about an HTML table

```
<table>
```

→ Defines a table

```
<tr>
```

→ Defines a row in a table

```
<td>Dogs</td> <td>Cats</td>
```

→ Defines a cell in a table

```
</tr>
```

```
</table>
```

Here is some information about furniture

```
<table>
```

```
<name>Australian Coffee Table</name>
```

```
<width>90</width>
```

```
<length>149</length>
```

```
</table>
```

What happens if we add these together in the one document?



- Namespace declarations are used to qualify names with **universal resource identifiers (URI's)**. A URI uniquely identifies a resource on the Web. The name consists of two parts
 - *namespace:local-name*
- This is achieved indirectly by using namespace declarations and associated user-specified prefixes

```
<... xmlns:tabular-info="http://www.tabularinfo.com"> <tabular-
```

```
info:table>
```

```
<tr>
```

```
<td>Dogs</td> <td>Cats</td>
```

```
</tr>
```

```
</tabular-info:table>
```

- **xmlns:tabular-info** attribute declares namespace with prefix tabular-info
- URI doesn't have to refer to a **real Web resource**

- The scope of a namespace declaration is
 - The **element** that contains the namespace declaration
 - All **its descendants** (i.e. nested within the element)
 - The declaration may be overridden by further nested namespace declarations
- Namespaces can be used to describe **both elements and attributes**. Elements without a namespace prefix are defined a default namespace.

```
<collection xmlns="http://www.tabularinfo.com"
xmlns:furniture="http://www.furniture.com">
```

```
<table>
<tr><td>Dogs</td> <td>Cats</td> </tr>
</table>
```

```
<furniture:table>
<furniture:name>Australian Coffee Table</furniture:name>
<furniture:width>90</furniture:width>
<furniture: length>149</furniture:length>
</furniture:table>
</collection>
```

- collection, first table, td and tr use the default (tabularinfo namespace)
- second table, name, width and length use the furniture namespace

```
<a:Envelope
  xmlns="http://default/"
  xmlns:a="http://urla"
  xmlns:b="http://urlb"
  xmlns:c="http://urlic"
  a:encodingStyle="http://schemas.xmlsoap.org/soap/encoding/">
  <a:Header xmlns="" xmlns:b="http://alturlb">
    <b:type>HelloWorld</b:type>
    <c:to xmlns:c="http://alturic">John Doe</c:to>
    <from fromType="name">Jane Seymour</from>
  </a:Header>
  <a:Body>
    <text xmlns="http://newdefault">Hello</text>
    <b:mood>Tired</b:mood>
    <c:day>Thursday</c:day>
    <month>March</month>
  </a:Body> <
/a:Envelope>
```

For each of the following, give its namespace URI: i) a:Envelope ii) a:Header iii) a:encodingStyle iv) b:type v) month vi) from vii) a:Body viii) text ix) b:mood

- We need to ensure the integrity of our data – define its expected structure and content
 - “A book element must have as children, a title, an ISBN and at least one author.”
 - “A title is a sequence of characters”, “An ISBN is ...”
- The format of the data can be specified by a **schema** and a document validated using schema checking software
 - Browsers use the HTML 5 Schema (see <!DOCTYPE html> at the start of an HTML document)
 - Schemas also used for other data formats
 - XML Schema (a W3C standard)
 - Large and complex, uses regular expression like rules
 - We will not look at the details in this subject



- An XML instance file is valid if it is consistent with a particular Schema
- Validation Tools
 - local XML editors (XMLWriter, Editix, Liquid XML ...)
 - online validators: <http://validator.w3.org/>
 - lxml (python library)
- **Note: an XML file can be well-formed and NOT valid**



- For HTML scraping, the BeautifulSoup library is good
- For XML, a good library is
 - <http://lxml.de/>
- Import the XML file into your program as a tree structure:

```
import xml.etree.ElementTree as ET
tree = ET.parse('yourfile.xml')
root = tree.getroot()
```

- Then loop through root with the various methods available:

```
for child in root:
    print child.tag, child.attrib
```



- Document Object Model (DOM)
 - Most useful way of parsing XML
 - Parsing calls load the document into a tree structure with different nodes that can be navigated by the program
- Simple API for XML (SAX)
 - Stream-based way of reading XML
 - Fast and efficient if you don't need random-access



- Further reading
 - Relational databases
 - Pages 403-409 of <http://i.stanford.edu/~ullman/focs/ch08.pdf>
 - XML
 - <http://www.tei-c.org/release/doc/tei-p5-doc/en/html/SG.html>