

# Advanced R Programming - Lecture 5

Krzysztof Bartoszek  
(slides by Leif Jonsson and Måns Magnusson)

Linköping University  
*krzysztof.bartoszek@liu.se*

18 September 2017

# Today

Input and output

Basic I/O

Cloud storage

web APIs: Lab

web scraping

Shiny

Relational Databases

# Questions since last time?

# Input and output



# Input and output



Format, localization and encoding..... hell!

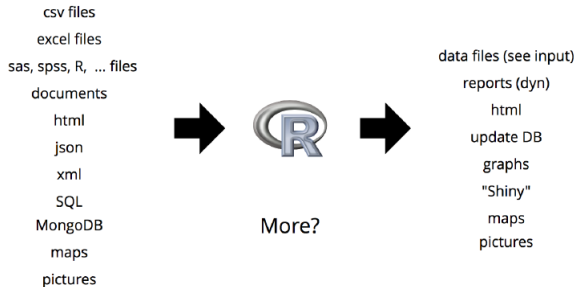
<http://www.joelonsoftware.com/articles/Unicode.html>  
The Absolute Minimum Every Software Developer Absolutely, Positively  
Must Know About Unicode and Character Sets (No Excuses!)

Unicode defines codes for **all (?)** characters—multiple encodings  
(for a given language only small fraction of characters used)

Content-Type tag for HTML

**BUT** e-mail, .txt, .csv

# "Formats"



# Localization



own Computer  
local network  
local database



Cloud Storage  
web pages  
web scraping  
web APIs  
remote database

Table: Local - Remote

## Files on your computer

```
# Input simple data
read.table()
read.csv()
read.csv2()

load()

# Output simple data
write.table()
write.csv()
write.csv2()

save()
```



## More complex formats

### **software/data**

Excel

SAS, SPSS, STATA, ...

XML

JSON (GeoJSON)

Documents

Maps

Images

### **package**

XLConnect

foreign

xml

rjsonio, RJSON

tm

sp

raster

Table: Format - R package

# Cloud storage



Table: Local - Remote

# Why?

Robust

Backups

Cloud computing

can be tricky in the beginning

**but**

# Why?

Robust

Backups

Cloud computing

can be tricky in the beginning

**but** how about safety?

**But** control on what is going on?

**BUT**

# Why?

Robust

Backups

Cloud computing

can be tricky in the beginning

**but** how about safety?

**But** control on what is going on?

**BUT** requires internet connection

# Localization

Arbitrary data



Structured data



# API Packages

<b>Remote</b>	<b>package</b>
General	downloader
GitHub	repmis, downloader
Dropbox	rdrop
Amazon	RAmazonS3
Google Docs	googlesheets

# web APIs

application program interface using http

"contract to 'get data' online"

more and more common

**examples:**

github

Riksdagen

Statistics Sweden



# RESTful

## Basic principles:

Data is returned (JSON / XML)

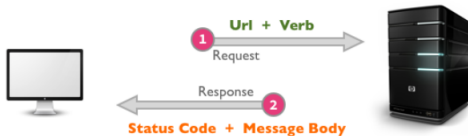
Each specific data has its own URI

Communication is based on HTTP verbs

# Hypertext Transfer Protocol (http)



# Hypertext Transfer Protocol (http)



# Verbs

<b>Verb</b>	<b>Description</b>
GET	Get "data" from server.
POST	Post "data" to server (to get something)
PUT	Update "data" on server
DELETE	Delete resource on server

# Status codes

<b>Code</b>	<b>Description</b>
1XX	Information from server
2XX	Yay! Gimme' data!
3XX	Redirections
4XX	You failed
5XX	Server failed

## Example REST API's

`http://www.linkoping.se/open/data/Luftkvalitet/  
Linköping Luftkvalitet API`

`https://developers.google.com/maps/documentation/geocoding/intro  
Google Map Geocode API`

# Common API formats

## **JavaScript Object Notation (JSON)**

Think of named lists in R

R Packages: RJSONIO, rjsonlite

## **Extensible Markup Language (XML)**

Older format (using nodes)

xpath

R Packages: XML

# JSON

```
{
  "firstName": "John",
  "lastName": "Smith",
  "age": 25,
  "address": {
    "streetAddress": "21_2nd_Street",
    "city": "New_York",
    "state": "NY",
    "postalCode": "10021"
  },
  "phoneNumber": [
    { "type": "home", "number": "212_555" },
    { "type": "fax", "number": "646_555" }
  ],
  "newSubscription": false,
  "companyName": null
}
```



# XML

```
<?xml version="1.0" encoding="utf-8"?>
<wikimedia>
<projects>
<project name="Wikipedia" launch="2001-01-05">
<editions>
<edition language="English">en.wikipedia.org</edition>
<edition language="German">de.wikipedia.org</edition>
<edition language="French">fr.wikipedia.org</edition>
<edition language="Polish">pl.wikipedia.org</edition>
<edition language="Spanish">es.wikipedia.org</edition>
</editions>
</project>
<project name="Wiktionary" launch="2002-12-12">
<editions>
<edition language="English">en.wiktionary.org</edition>
<edition language="French">fr.wiktionary.org</edition>
<edition language="Vietnamese">vi.wiktionary.org</edition>
<edition language="Turkish">tr.wiktionary.org</edition>
<edition language="Spanish">es.wiktionary.org</edition>
</editions>
</project>
</projects>
</wikimedia>
```

# web scraping

Unstructured http(s) data

Often HTML format

Spiders / scraping / web crawlers

Basics behind search engines

# HTML

```
<!DOCTYPE html>
<html>
  <head>
    <title>This is a title</title>
  </head>
  <body>
    <p>Hello world!</p>
  </body>
</html>
```

# (har)rvest

## JavaScript Object Notation (JSON)

Simplify spider activity

Download data

Parse data

Follow links

Fill out forms

Store crawling history

# Difficulties and bad spiders

Scraping is fragile!

Difficulties and bad spiders

[www.domain.se/robot.txt](http://www.domain.se/robot.txt)

Politeness

robot traps

javascript

delays

# Shiny?

Interactive dashboards made easy

online or local

R as "backend"

# Shiny?

<https://www.rstudio.com/products/shiny/shiny-user-showcase/>  
Shiny Examples

# How it works

Application

Reactive

modify using HTML

```
MyAppName/server.R
```

```
MyAppName/ui.R
```

server.R define working directory



# Shiny Example

```
library(shiny)
# Examples with code
runExample("01_hello")
runExample("03_reactivity")
```

# Publish Shiny



locally  
zip-file in cloud  
github (see `runGithub()` )

# Publish Shiny



locally  
zip-file in cloud  
github (see `runGithub()` )



your own server  
shinyapps.io

# Relational Databases

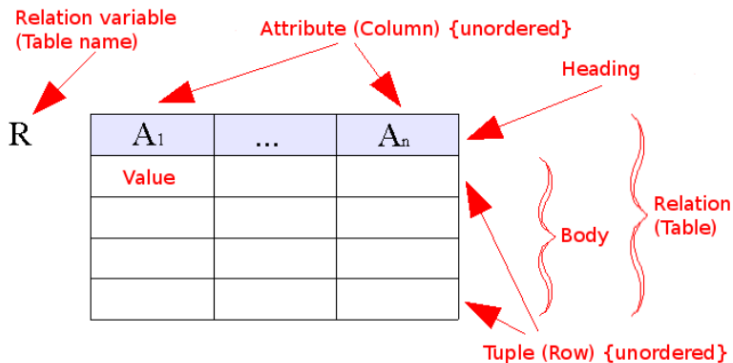
Structured database in tables

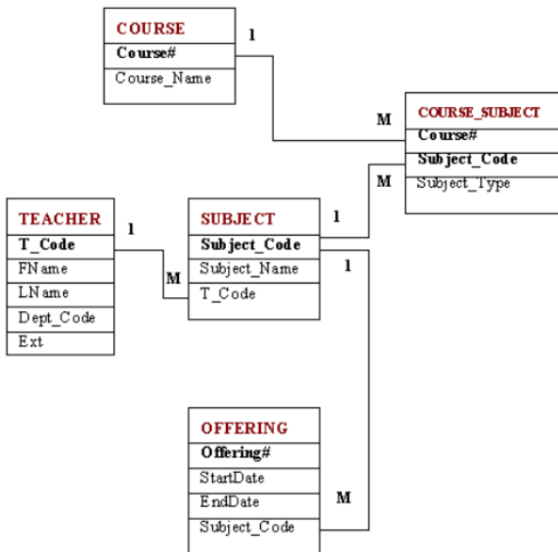
local or online

query language for I/O

effective for big data

difficult to design





# A good database

Can be difficult to design ?

# A good database

Can be difficult to design ?

No duplicates

No redundancies

Easy to update

"Normal forms"



# A good database

Can be difficult to design ?

No duplicates

No redundancies

Easy to update

"Normal forms"

Easy to query

# Using databases from R

<b>Database system</b>	<b>R package</b>
ODBC (Microsoft Access)	RODBC
PostgreSQL	RPostgresql
Oracle	ROracle
MySQL	RMySQL
MongoDB	rmongodb

**Table:** Database - R package

The End... for today.  
Questions?  
See you next time!