

UNIVERSITÀ CATTOLICA SACRO CUORE

STATISTICAL AND ACTUARIAL SCIENCES

MJ: DATA BUSINESS ANALYTICS

Spatial Machine Learning modelling: End-to-End web App solution

Author:

Niccolò SALVINI

Supervisor:

Dr. Marco DELLAVEDOVA

Assistant Supervisor:

Dr. Vincenzo NARDELLI

AY 2019 / 2020



Spatial Machine Learning modelling: End-to-End web app solution

Niccolò Salvini¹

date: Last compiled on 17 settembre, 2020

¹<https://niccolosalvini.netlify.app/>

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 2 |
| 2 | Scraping | 6 |
| 2.1 | What is Scraping | 6 |
| 2.2 | User agents, Proxies, Handlers | 8 |
| 2.3 | How they are designed with <code>rvest</code> | 14 |
| 2.4 | What are the advantages of workflow | 14 |
| 2.5 | Legal challenges | 14 |
| 3 | Infrastructure | 16 |
| 3.1 | Scheduler | 19 |
| 3.2 | Docker Container | 21 |
| 3.3 | API | 28 |
| 3.4 | What is an API | 28 |
| 4 | Methodologies | 33 |
| 5 | Applications | 38 |
| 5.1 | Example one | 38 |
| 5.2 | Example two | 38 |
| 6 | Final Words | 39 |

Chapter 1

Introduction

Main themes:

- (open data)
- research question
- milan market real estate
- latest improvements in the subject matter
- why both bayesian and non bayes methods
-

We are living in the big data era, so we could be brought to think that everything is a “one click” distant from us. Well, this is not totally true, moreover in some places this is truer. The main issue can be addressed to the lack of open data and the lack of relative infrastructure. This settings characterizes slow old economies and unfortunately Italy is one of them. Economies, and citizens on a later step, can largely benefit from public data and its usage. Some people in addition are in favor of the position that all data should be open. Since I am living in italy and my (Lovelace et al., 2019) goal is to analyse market

The importance of data indeed justifies its accumulation and according to the latest reports is surpassing gold, despite these periods of uncertainty. The expression data is the new oil has never been so appropriate in these times. On the other hand is not for sure easy to assign a price amount to data due to its untangible nature. the most straightforward and liberal approach could lead us to think that the price data should be exchanged the piceThe value attributed to data is not for sure selrmarkdown::pandoc_available("1.2") f explanatory. It really depends on two major metrics: the usage that can be done through (with respect of the state of the art technology) it and the functionality with respect of other existing data. some data can be strategically important given the fact that someone already possess the complementary and can attribute some sort of competitive advantage. On the other hand as already been highlighted it really depends on the existing technology stack. Some data can be very useful but too costly either to process or to store.

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k} \quad (1.1)$$

During an interesting conversation with some friends we had a discussion on how data should be treated: as a sort of currency or a sort of commodity (raw material). some people may say that the inner functioning is pretty much as a commodity. It gains value by its specialized usage and treatment. Sometimes a collection of data can represent the complementary part of a more general dataset that can not be used otherwise, in this analogy case a semi-finished product. commodities sometimes are calmed, so that their prices are fixed to a certain amount, so it is for data, the

Durante una conversazione con alcuni amici¹

La ricerca che ho inteso fare sul mercato degli affitti a Milano mi ha aperto le porte a comprendere come poco digitalizzata e all'avanguardia sia la nostra amata penisola. L'indisposizione ai dati aperti, coperta da un sottile velo di ipocrisia chiamato privacy (ma quale?), ha reso non solo impossibile reperire

¹footnote a caso

i dati geospaziali tramite API di alcune aree dell'italia, ma ha reso necessario che costruissi delle funzioni che li estressero, appoggiandomi a cavilli. La questione è legale e relativamente complessa, e di certo la tesi non si indirizza a questi problemi, ma i dati che sono stato capace di scoprire e di farlo nella assoluta legalità si appoggiano ad una mancanza di autorizzazioni al trattamento che Immobiliare.it ha nel suo sito. Un altro esempio di ritardo tecnologico riguarda l'assenza dei dati di elevazione su alcuni territori italiani. Se per esempio immobiliare, come ha fatto un altro grosso player sul mercato, avesse apposto una checkbox obbligatoria da contrassegnare con relativi termini e trattamento dei dati io non avrei potuto accedere ai dati. La situazione al di fuori dell'italia è abbastanza uniforme, eccetto qualche paese noto come Germania e Francia, e meno noto come la Polonia, con tutte altre piattaforme e regole di trattamento dati. La domanda quindi sorge spontanea, perchè i dati degli italiani e degli europei sono meno accessibili dei dati degli americani? Mentre in America è sufficiente richiamare un API con latitudine e longitudine della quadrettatura di terra necessaria per ottenere i dati² di elevazione (.tif), in Italia l'unica soluzione è pagare google che tramite le sue private API è in grado di venderceli dietro autenticazione. La risposta al di là dei confini della legge presumibilmente risiede in un congiunturale ritardo di infrastrutture tecnologiche condivise e di indirizzo comune europeo sulla questione. L'esigenza di dati aperti nasce per la risoluzione di problemi comuni a tutti, i dati sanitari hanno la missione di tentare risolvere problemi di natura sanitaria, i dati economici auspicabilmente curano problemi o asimmetrie di un mercato. Il mercato degli affitti a Milano gode di sempiterna gloria e ha visto la crescita degli affitti e dei prezzi degli immobili di paripasso al punto che una bolla è stata presunta. Diversi fattori hanno reso tale il fenomeno e diverse opinioni si sono spese sul tema. Alcuni pensano che dopo Expo la città abbia goduto di una spinta economia e innovativa che l'ha resa un'isola felice in mezzo ad un'italia che affanna. Altri ritengono che Milano goda di ottime infrastrutture, ma che la sua notorietà ed il suo appeal si sia sostituito a tutto

²<https://it.wikipedia.org/wiki/Dato>

quello che manca nelle altre città, ma che in Milano appare. La mia opinione è che sia una media di questi due pareri. Un altro fattore è importante nella descrizione del fenomeno: l'asimmetria di informazione tra chi cerca casa a Milano venendo da fuori e colui che affitta. Tale asimmetria viene ancora più esasperata al crescere della fretta che l'entrante ha nel trovare la locazione opportuna. La scelta diventa in molti casi antieconomica, nello specifico la domanda si genuflette all'offerta e accetta le svattaggiose condizioni proposte. Infatti quello che appare certo è che i prezzi degli affitti se comparati ai salari per posizioni junior e di stage è falsato. Proprio qui nasce l'esigenza di approfondire il perchè e fornire all'utente finale (un potenziale studente, un futuro lavoratore etc.) uno strumento che gli permetta di capire il prezzo stimato tramite predizione spaziale date le coordinate geografiche e gli attributi dell'appartamento e contestualmente fornire un mezzo di comparazione per altri immobili nelle vicinanze. Dall'altro lato dia un'idea chiara a chi vuole dare in affitto l'immobile, un prezzo rappresentativo, che ha fondamento nel modello utilizzato e nelle assunzioni che lo stesso modello impone alla realtà. Questo fa sì che da entrambi i lati ci sia trasparenza e che eventuali maggiorazioni di prezzo richiesto rispetto al sopradetto modello vengano penalizzate in favore di sconti applicati su altri immobili. Auspicabilmente i prezzi già gonfi si smusseranno in tutta la regione spaziale considerata, adattandosi alla domanda piuttosto che al capriccio dell'offerta.

Chapter 2

Scraping

2.1 What is Scraping

Lo web scraping è una tecnica di estrazione dei dati da pagine internet statiche o dinamiche in maniera automatica e simultanea (Wikipedia, 2020). L'impossibilità di reperire dati aperti aggiornati riguardo l'affitto sul mercato italiano mi ha spinto a sviluppare sofisticate tecniche di estrazione di dati orientate ad alleggerire lo sforzo e aumentare la velocità di reperimento: da una parte nel preprocessing del dataset, nella successiva del frangente del modelling, per finire con la reattività di risposta dell'applicazione. Le informazioni sui siti appaiono spesso ordinate e semplici, tuttavia ogni sito web ha una propria architettura e un proprio linguaggio. Per architettura intendo struttura gerarchica secondo cui è organizzato un sito internet: una semplificazione della struttura di un sito web può essere un insieme di cartelle innestate una dentro l'altra collegate tra loro da riferimenti tramite l'url. la natura gerarchica della struttura prevede che si usi un linguaggio che fa propria questa caratteristica, HTML è il preferito. L'html si organizza in nodi ed angoli, esattamente come un grafo; che aggiunta la componente gerarchica fa sì che questo sia un albero. Difatti spesso ci si riferisce alla struttura delle pagine web come html tree. Ogni elemento nella pagina ha un suo preciso posto nel codice sorgente della stessa e ha un preciso valore o più valori.

Possiamo immaginare ogni nodo della pagina come una lista di valori che è collegata ad un nodo precedente detto padre da una struttura gerarchica superiore, ed eventualmente ad un nodo successivo detto figlio. Pertanto tutte le informazioni che giacciono sotto al nodo padre sono parenti del nodo padre e sono direttamente collegate (directed nel senso dell'interpretazione), parallelamente ci saranno altri nodi padre che saranno adiacenti al nodo padre, i quali avranno nodi figli e così via. La complessità della pagina e del codice è tanto maggiore quanto il livello dell'albero aumenta, tanto più l'albero è folto tanto più sarà difficile individuare il ramo o la foglia che ci interessa. Ragionevolmente accade lo stesso per la funzione di scraping e il tempo di scraping. Html organizza i contenuti e le relazioni tra loro, il css (Cascading Style Sheets) invece si occupa dello stile e della formattazione degli stessi. il css è uno strumento molto potente in mano ad uno scraper perchè permette di recuperare informazioni simili tra loro ma che occupano nodi con posizione gerarchica diversa all'interno della pagina. Pertanto una volta letto l'html della pagina sarà necessario recuperare la query css per raccogliere tutti gli elementi di interesse tramite la funzione di scraping. Successivamente occorre notare che l'encoding da html a stringa di testo non è quasi mai lineare, spesso occorre riformattare, cancellare spazi, convertire la natura dell'oggetto estratto etc. Il successivo elemento di complessità incontrato durante questa prima fase è stato interfacciarsi con un server attento alle richieste GET degli utenti. I dati viaggiano in pacchetti da un server che ospita un sito internet al nostro laptop. tutte le volte che cerchiamo di accedere ad un sito stiamo mandando una richiesta di ricezione di pacchetti dati ad un server in qualche luogo remoto del mondo. Quando bussiamo alla porta del server se non siamo sospetti e superiamo i criteri autostabiliti dal server questo risponde, e lo fa con un numero che spazia da 200 a 500, due esempi: 200 se la risposta è positiva, 404 se la risposta è negativa. I criteri secondo cui gli utenti sono classificati secondo utente normale o utente sospetto (aka bot) sono sintetizzati in un documento di testo chiamato robot.txt. Questo file di testo raccoglie tra le altre due informazioni principali il delay time, cioè il tempo

preferito dal server che deve intercorrere tra una richiesta dati e la successiva e quale utente è autorizzato ad accedere. Ogni utente possiede un indirizzo IP che nelle richieste a server si codifica in user agent, cioè una stringa di testo dove vengono raccolte le informazioni significative circa il dispositivo da cui provengono le richieste, un esempio:

‘Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/54.0.2840.71 Safari/537.36’,

dove ogni segmento della stringa rispecchia una caratteristica del laptop del richiedente, Chrome/54.0.2840.71 è la versione del browser chrome da cui proviene la richiesta Safari/537.36’, è il motore di ricerca etc.

2.2 User agents, Proxies, Handlers

Everytime a user enters a website what he is really doing is sending an HTTP request to the website server with some information packed. This can be easily thought as a generic person A that rings the door’s bell of person B’s house. A comes to the B door with its personal information, its name, surname, where he lives etc. At this point B may either answer to A requests by opening the door and let him enter given the set of information he has, or it may not since B is not sure of the real intentions of A. This typical everyday situation in nothing more what happens billions of times on the internet everyday, the user (in the example above A) is interacting with a server website (part B) sending packets of information. If a server does not trust the information provided by the user, if the requests are too many, if the requests seems to be scheduled due to fixed sleeping time, a server can block the requests. In certain cases it can even forbid the user to be on the website. The language the two parties talks are coded in numbers that ranges from 100 to 511, each of which ha its own significance. A popular case of this type of interaction occurs when users are not connected to internet so the server responds 404, page not found. Servers are built with a immune-system like software that raises barriers and block

users to prevent dossing or other illegal practises.

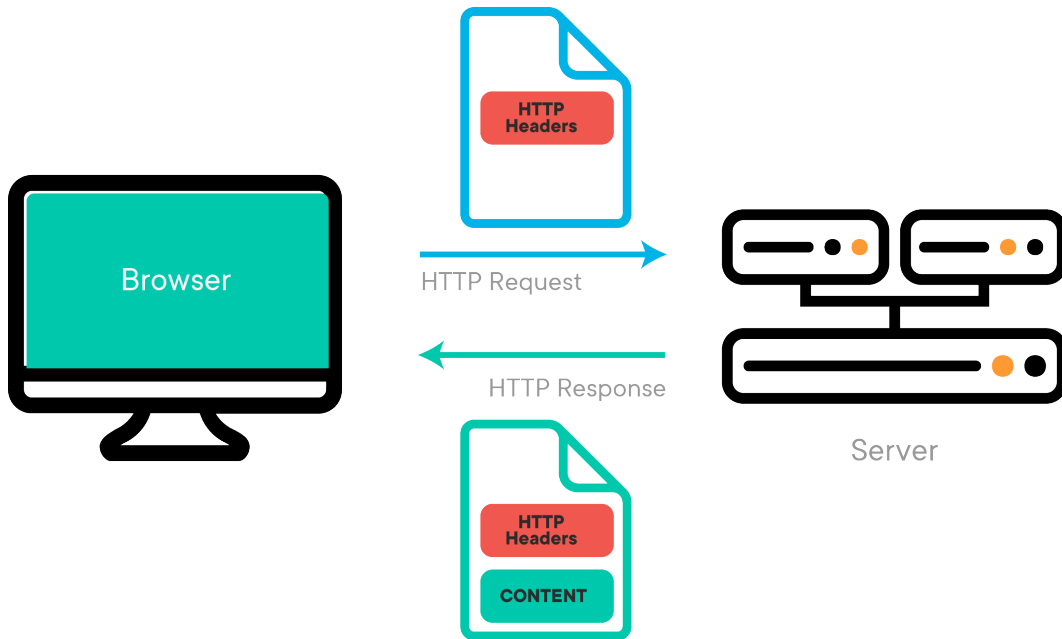


Figure 2.1: How Web Works

This procedure is a daily issue to people that are trying to collect information from websites. Google does it everyday with its spider crawlers, which are very sophisticated bots that performs scarping over a enormous range of websites. This challenge can be addressed in multiple ways, there are some specific Python packages that overcome this issue. There are also certain types of scraping as the Selenium web driver automation that simulates browser automation. Selenium allows the user not to be easily detected by the server immune system and peaceful. In here precautions have not been taken lightly, and a simple but effective approach is proposed.

2.2.1 User agents Spoofing

A user agent (who, 2020) is a string of characters in each browser that serves as an identification agent. The user agent permits the web server to be able to identify the user operating system and the browser. Then, the web server uses the exchanged information to determine what content should be presented to particular operating systems and web browsers on a series of devices. The user

agent string contains the user application or software, the operating system (and their versions), the web client, the web client's version, and the web engine responsible for the content display (such as AppleWebKit). The user agent string is sent in form of a HTTP request header. Since the User Agents acts as middle man between the client request and the server response what it would be better doing is to actively faking it so that each time a web browser presents himself to a web server it has a different specifications, different web client, different operating system and so on.

The simple approach followed was building a vector of samples of different existing and updated User Agents (UA). Then whenever a request from a browser is served to a web server 1 random string is drawn from the user agents pool. So each time the user is sending the requests it appears to have a different “identity”. Below the user agents rotating pool:

```
set.seed(27)
agents = c('Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_12_1) AppleWebKit/602.2.14 Safari/602.2.14',
           'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_12_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_6) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Windows NT 6.1; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36',
           'Mozilla/5.0 (Windows NT 10.0; WOW64; rv:50.0) Gecko/20100101 Firefox/50.0')

agents[sample(1)]
```

```
## [1] "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2866.90 Safari/537.36"
```

An improvement to this could be using also rotating proxies. A proxy server acts as a gateway between the user and the server. It's an intermediary server

himself, separating end clients from the websites they are browsing. Proxy servers provide varying layers of functionality, security, and privacy are some of the examples. While the user is exploiting a proxy server, internet traffic flows through the proxy server on its way to the server you requested. The request then comes back through that same proxy server and then the proxy server forwards the data received from the website to you. Many proxy servers are offered in a paid version so in this case since security barriers of the target website are not high they will not be implemented. It has to be mentioned that many online services are providing free proxies but the turnaround of this solutions are many, two of them are: - Proxies to be free are widely shared among people, so as long as someone has used them for illegal purposes the user is inheriting their mistakes when caught. - Some of those proxies, pretty all the ones coming from low ranked websites, are tracked so there might be a user privacy violation issue.

2.2.2 Handlers

During the scraping many things could be going wrong. Starting from the things that have been previously explained at the chapter start (URL structure changes, data are moved in different location so that css query goes empty...), ending with the ones that have just been said a few lines ago. Handlers and trycatch error workarounds are explicitly built in this sense. The continuous testing of the scraping functioning while developing has required the maintainer to track where the error occurs. A few numbers: the “agglomerative” function `get.data.catsing()` triggers more than 36 scrapping functions that are going to catch 36 different data pieces. If one of them went missing then the other one would be missing too. Then when row-data is binded together one entry column might not exists making the process fail.

Then the solution to that is to call inside the agglomerative function as much as trycatch as many scrapping functions are involved. The trycatch can leverage the gap by introducing a specified quantity and alerting that something

went wrong. On top of that many other handlers are called throughout the procedure:

- `get_ua()` verifies that the user agent coming from the session request is not the default one

```
get_ua = function(sess) {  
  stopifnot(is.session(sess))  
  stopifnot(is_url(sess$url))  
  ua = sess$response$request$options$useragent  
  return(ua)  
}
```

- `is_url()` verifies that the url input needed has the canonic form. This is done by a REGEX query.

```
is_url = function(url){
  re = "^((?:http(?:s)?|ftp)://)(?:\\S+(?:\\S*)?@)?(?:[a-z0-9\\u00a1-\\u00ff-"]
  grepl(re, url)
}
```

- `.get_delay()` checks through the robotxt file if a delay between each request is welcomed.

```
.get_delay = function(domain) {

  message(sprintf("Refreshing robots.txt data for %s...", domain))

  cd_tmp = robotstxt::robotstxt(domain)$crawl_delay

  if (length(cd_tmp) > 0) {
    star = dplyr::filter(cd_tmp, useragent=="*")
  }
}
```

```
    if (nrow(star) == 0) star = cd_tmp[1,]
    as.numeric(star$value[1])
  } else {
    10L
  }
}

get_delay = memoise::memoise(.get_delay)
```

- `checkpermission()` checks through the `robot.txt` file if the dom allows the user to scape data.

```
dominio = "immobiliare.it"

checkpermission = function(dom) {

  robot = robotstxt(domain = dom)
  vd = robot$check()[1]
  if (vd) {
    cat("\nrobot.txt dice: va bene, puoi!")
  } else {cat("\nrobot.txt dice: non puoi, smettila")}
}
```

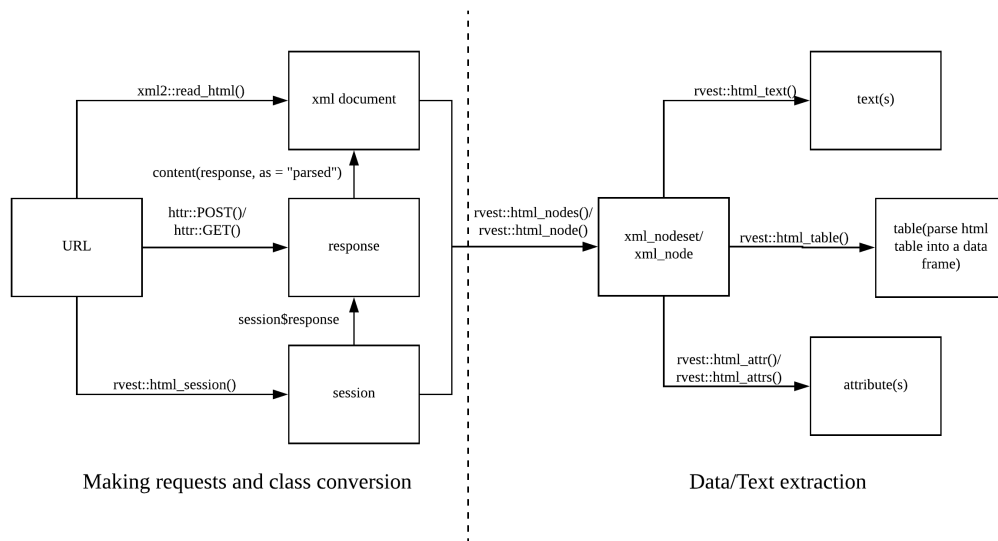


Figure 2.2: functional structure

2.3 How they are designed with rvest

2.3.1 From Generic and Specific structure

2.3.2 parallel computing

2.3.3 foreach package and runtime benchmark with future

2.4 What are the advantages of workflow

2.5 Legal challenges

“Data that are online and public are always available” is never a good answer to the question “Can I use that data to my scope”. Immobiliare.it¹ is not providing any source of data from its own database neither it is planning to do so in the future.

¹<https://www.immobiliare.it/>

- Copia e incolla manuale
- Web scraper
- HTML parsing
- Analisi con Visione computerizzata
- DOM parsing
- Riconoscimento dell'annotazione semantica
- Aggregazione verticale
- Text pattern matching

funzioni di scraping.

Chapter 3

Infrastructure

In order to provide a fast, portable and integrated product to the end user It has been designed a quite straightforward software architecture. We have already seen the scraping functions and how they are built around the concept of easy flexibility and debugging. This is due to the fact that they should extract something that is dynamic, it is not sure that it will be as the day before. The data we are trying to grab might have been moved somewhere else throughout the website. Or it might have placed extra expression inside the node we are inspecting (“\$” sign following the monthly rental price) . A very often occurring example regards the way information concerning the house are represented in the website. Considering the september 2018 january 2020 time span the design of the website has changed a vast number of times. Since both the design and the scrapping functions relies on the HTML skelethon and CSS queries. As soon as something changes in the website the other files needs to be readjusted to be consistent with the content and so back and forth. The debugging handlers nested in the functions helps the maintainer to grasp what it is not working properly where the error occur. The same inner philosophy has been applied to the software architecture chosen for this project. First of all the wide range of open source solutions (back-end and front-end) and documentation on this has made many analyst and data scientist almost full stack developer. This was also due to the fact that RStudio has set very well

oriented guidelines spending a lot of effort giving its users an easy, integrated and interconnected environment. By that it is meant that recently the RStudio community has developed, on top of many different others, an entire package dedicated to REST APIs (Plumber (Trestle Technology, LLC, 2018)). Moreover developers in RStudio and its contributors have created an entire new paradigm called Shiny (Chang et al., 2020), a popular web app development package, that forces the user to have front-end and back-end technologies tied up in the same IDE (RStudio) and with a unique language to deal with. The front end file (for simplicity named UI.R) contains the UI's layout and the style and also other javascripts components. On the other hand the server file (named after server.R) absorbs the back-end, under the hood, code and makes the UI interact and respond to the user. This comes at a cost of flexibility and customization since Shiny could not easily handle too many embellishments (even though potentially can). Nevertheless a unique environment makes integration with other technologies easier and most of all introduces the analyst to a full stack approach. Many open source projects are gravitating around the Shiny framework with the aim to extend its capabilities. One example is a newly created package called reactR (Inc et al., 2020) that allows user to implement the power of React.js into the shiny UI front end. All of this is possible, once again, by the R community but a greater contribution come from digging up the right path along which everything by open source comes natural. (parallelo con la vigna e l'albero che la sostiene e indirizza) The carrier idea for this project is to have parallelized scraping functions called daily by a scheduler producing and subsequently storing a .csv file in a MySQL /cartoDB database. They are all thought to be containerized in a Linux (Ubuntu distr) docker container hosted in a AWS EC2 server. Then in a second container a Shiny app is placed, this one pipes in data from the former infrastructure and apply the statistical model stored (by an API call) in its server.R part.

The main technologies implied are:

- Scheduler cron job

- Docker containers
- Shiny
- Plumber REST API
- AWS (Amazon Web Services) EC2
- CartoDB

On top of that even each single part of this thesis has been made stand alone and can be easily accessed and modified through this link¹. The pdf (theis) version of the gitbook can be obtained by clicking the download button that can be seen in figure below. A Latex engine (Xelatex) wrapped into the website compiles a sequence of Markdown documents converting them into .html (the book's chapters) which are formatted by rules grouped in a .yaml file. All the documents are pushed to a Github repository with git. By a simple trick, since all the files are static html, they can be displayed through GH pages as it is a website. All of this has been possible thanks to Bookdown (Xie, 2020) once again a R well documented package (Xie, 2016) to build interactive books along with RMarkdown (Allaire et al., 2020).

An empirical observation of immobiliare.it has suggested that houses rents advertisement are continuously added and then removed during the day. Fresh data is needed to have updated analysis since the scope in here is to offer realtime considerations. Something should be automated periodically in order to address the issue. Moreover, as rule of thumb, a daily data extraction might be a good option for some reasons. It can intercept price variations with a relatively small time lag, It can also display some sort of pattern in time that would help the reader/user to select the perfect choice. As a consequence a daily .csv file is generated and directly collected into a Db folder arranged by time The solutiond proposed takes care of the issue by making the scraping script generating the .csv be executed by a scheduler.

¹<https://niccolosalvini.github.io/Thesis/>

3.1 Scheduler

A Scheduler in a process is a component on a OS that allows the computer to decide which activity is going to be executed. In the context of multi-programming it is thought as a tool to keep CPU occupied as much as possible. As an example it can trigger a process while some other is still waiting to finish. There are many type of scheduler and they are based on the frequency of times they are executed considering a certain closed time neighbor.

- Short term scheduler: it can trigger and queue the “ready to go” tasks
 - with pre-emption
 - without pre-emption

The ST scheduler selects the process and It gains control of the CPU by the dispatcher. OIn this context we can define latency as the time needed to stop a process and to start a new one.

- Medium term scheduler
- Long term scheduler

for some other useful but beyond the scope information, such as the scheduling algorithm the reader can refer to (Wikiversità, 2020).

The scheduler in this context cosists in a .sh (shell file, sort of text file) composed by a set of instructions that are being executed by the computer on daily basis. This file has to be in the same WD (working directory) of the project in order to make it working. Some common issues can occur when new files coming after the execution of the scheduled main script are generated, but the path isnt explicitly specified. This can lead to the partial or incomplete generation of the file since the shell file is executed within the folder but is triggered by some other location on the computer. Each OS has its own scheduler and syntax to call it. Since we are interested in Ubuntu machines the scheduler

is said to be a cron job. Later it will be clear why Ubuntu is the option to pursue.

va parafrasato

3.1.1 Cron Jobs

The software utility cron also known as cron job is a time-based job scheduler in Unix-like computer operating systems. Users that set up and maintain software environments use cron to schedule jobs (commands or shell scripts) to run periodically at fixed times, dates, or intervals. It typically automates system maintenance or administration—though its general-purpose nature makes it useful for things like downloading files from the Internet and downloading email at regular intervals.

The actions of cron are driven by a crontab (cron table) file, a configuration file that specifies shell commands to run periodically on a given schedule. The crontab files are stored where the lists of jobs and other instructions to the cron daemon are kept. Users can have their own individual crontab files and often there is a system-wide crontab file (usually in `/etc` or a subdirectory of `/etc`) that only system administrators can edit.

Each line of a crontab file represents a job, and looks like this:

```
# _____ minute (0 - 59)
# _____ hour (0 - 23)
# _____ day of the month (1 - 31)
# _____ month (1 - 12)
# _____ day of the week (0 - 6) (Sunday to Saturday;
#                                     7 is also Sunday on some systems)
#
# * * * * * <command to execute>
```

Figure 3.1: crontab

Each line of a crontab file represents a job. This example runs a shell program called `scheduler.sh` at 23:45 (11:45 PM) every Saturday.

```
45 23 * * 6 /home/oracle/scripts/scheduler.sh
```

Some rather unusual scheduling definitions and syntax for cronjobs can be found in this reference (Wikipedia contributors, 2020)

The cron job applied to the script needs to be ran at 11:30 PM everyday. It has that forms: —> qui immaginare

va parafrasato

For now the computational power comes from the machine on which the system is installed. A smarter solution takes into consideration that the former infrastructure has its own limits. Major limits comprehend run time since at the same moment the machine runs locally both the scraping functions and the app computations. This to a certain extent might fit for personal use but as data increases all the system risks to fail. It is also totally local so the analysis can not be shared with anyone. This problem can be addressed with a technology that has seen a huge growth in its usage in the last few years: Docker containers.

3.2 Docker Container

from docker In 2013, Docker introduced what would become the industry standard for containers. A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another. A Docker container image is a lightweight, standalone, executable package of software that includes everything needed to run an application: code, runtime, system tools, system libraries and settings.

3.2.1 What is Docker?

Container images become containers at runtime and in the case of Docker containers - images become containers when they run on Docker Engine. Available for both Linux and Windows-based applications, containerized software will al-

ways run the same, regardless of the infrastructure. Containers isolate software from its environment and ensure that it works uniformly despite differences for instance between development and staging.

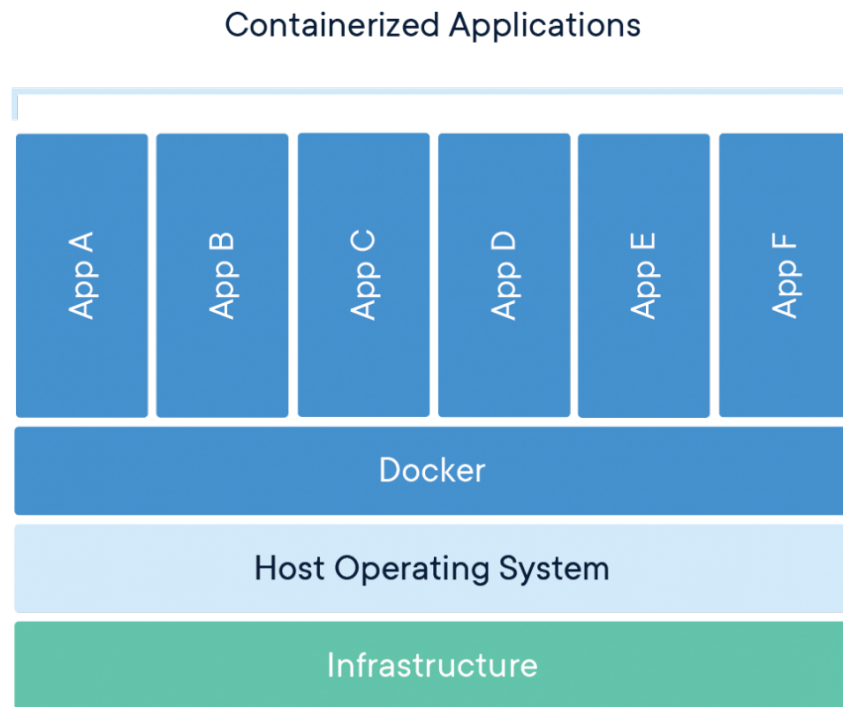


Figure 3.2: docker example

Docker leveraged existing computing concepts around containers and specifically in the Linux world. Docker's technology is unique because it focuses on the requirements of developers and systems operators to separate application dependencies from infrastructure.

A question might come up about why a Virtual Machine could not be a preferable container for our specified task. Well, Containers and virtual machines have similar resource isolation and allocation benefits, but function differently because containers virtualize the operating system instead of hardware. Containers are more portable and efficient.

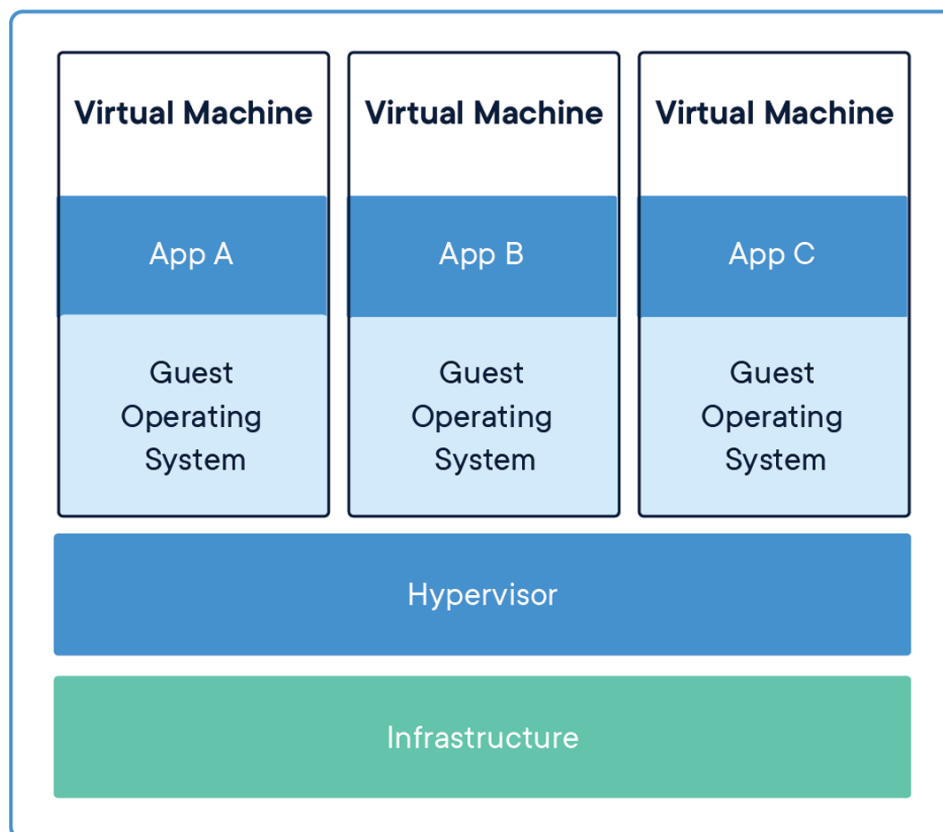


Figure 3.3: docker container vs VM

from docker

3.2.2 What are the main advantages of using Docker

va parafrasato from Matt Dancho

Indeed, the popular employment-related search engine, released an article this past Tuesday showing changing trends from 2015 to 2019 in “Technology-Related Job Postings”. We can see a number of changes in key technologies - One that we are particularly interested in is the 4000% increase in Docker.

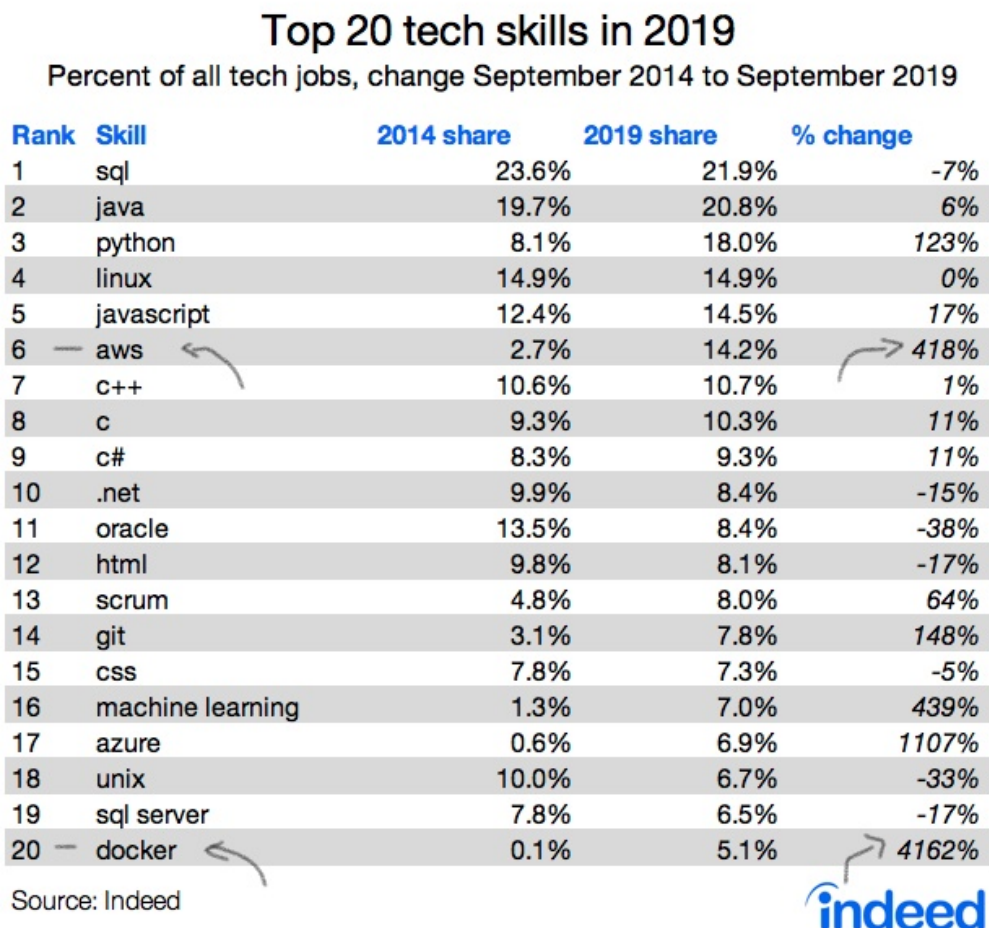


Figure 3.4: docker-stats

The landscape of Data Science is changing (was previously an Economist at the Indeed Hiring Lab, 2020) from reporting to application building:

In 2015 - Businesses need reports to make better decisions
In 2020 - Businesses need apps to empower better decision making at all levels of the organization
This transition is challenging the Data Scientist to learn new technologies to stay relevant...

As a matter of fact, it is no longer sufficient to just know machine learning algorithms. Future data workers need to know how to put machine learning into production as quickly as possible to meet the business needs. This can be done either integrating existing obsolete/old technologies with the new ones, or build a solid, portable and scalable infrastructure to better the processes. In order to do so, It is strictly needed to learn from programmers the basics of Software Engineering. The extremely good news is that no data scientist whatsoever needs to be a perfect software engineer, but at least he has to know how to integrate already built technologies with its work. This truly can help in the quest to unleash data science at scale and unlock real business value. The way Docker does that are many (red, 2020):

- *Rapid application deployment* : containers include the minimal run time requirements of the application, reducing their size and allowing them to be deployed quickly.
- *Portability across machines* :an application and all its dependencies can be bundled into a single container that is independent from the host version of Linux kernel, platform distribution, or deployment model. This container can be transfered to another machine that runs Docker, and executed there without compatibility issues.
- *Version control and component reuse* : you can track successive versions of a container, inspect differences, or roll-back to previous versions. Containers reuse components from the preceding layers, which makes them noticeably lightweight.
- *Sharing* : you can use a remote repository to share your container with others. It is also possible to configure a private repository hosted on Docker Hub.

- *Lightweight footprint and minimal overhead* : Docker images are typically very small, which facilitates rapid delivery and reduces the time to deploy new application containers.
- *Simplified maintenance* : Docker reduces effort and risk of problems with application dependencies.

The way all of this is possible is a dockerfile that determines the instruction that docker has to perform to abstract the environment.

3.2.3 Dockerfile

Docker can build images automatically by interpreting the instructions from a Dockerfile. A Dockerfile can be thought as a written recipe to cook a specific cake, with all the ingredients described in a piece of a paper using a generic oven and adding layer of preparation after layer. A Dockerfile is a text format document that contains all the commands/rules a generic user could call on the CLI to assemble an image. Executing the command `docker build` from shell the user can trigger the image building. That executes several command-line instructions in chronological succession of steps. The Dockerfile used to trigger the build of the docker image has this following set of instructions:

- `FROM rocker/r-ver:4.0.0` : the command imports an image already written by the rocker team (authored contributors for the R docker project) that contains the base-R version 4.0.0. Recently with the 4.0 version the RStudio team has created a repository management server for its packages that organizes and centralizes R packages (offline access and checkpoints). This will shorten the installation time and secure packages since they all can be freezed into a version that make the whole system works.
- `RUN R -e "install.packages(c('plumber','tibble','...',dependencies=TRUE))` : the command install all the packages required to execute the files (R



Figure 3.5: dockerfile

files) containerized for the scraping. Since all the packages have their dependencies the option `dependencies=TRUE` is needed.

- `EXPOSE 8000` : the command instructs Docker that the container listens on the specified network ports 8000 at runtime. It is possible to specify whether the port exposed listens on UDP or TCP, the default is TCP (this part needs a previous set up previous installing, for further online documentation It is recommended (doc, 2020))
- `ENTRYPOINT ["Rscript", "main.R"]` : the command tells docker to execute the Rscript extension file main.R within the container that triggers the API building/the generation of the .csv file.

va parafrasato from Matt Dancho

An alternative and very used approach could be wrapping all the scraping function into an API and then send a `GET` request to the API endpoint needed.

3.3 API

va parafrasato

The scraping functions, according to how the author as structured them, are able to produce two .csv extension (if the boolean option `write` is set = `TRUE`) files. As already clarified in the previous ssection some point should be joined by a primary key. But for the sake of In order to give the possibility to have a daily updated saptial analysis on data we need to continously have fresh data. In the website data come and go, as products in a marketplace, so the main idea is to have something that catches the new added and deletes what it is already taken. Nowadays we have many open source, nearly cost free, techonologies that allow us to have corporate grade applications that can be orizontally scaled at need. Most of them come with great docuemntation and ready to use examples that flatten the learning curve. The first choice that has to be made is: either to provide a .csv file day by day with all the data to feed the application, or we exploit some portable and fast solutions as API.

va parafrasato

3.4 What is an API

API is a set of definitions and protocols for building and integrating application software. API stands for application programming interface.

APIs let a product or a service communicate with other products and services without having to know how they're implemented. This can simplify app development, saving time and money. When you're designing new tools and products—or managing existing ones—APIs give flexibility; simplify design, administration, and use; and provide opportunities for innovation. APIs are sometimes thought of as contracts, with documentation that represents an agreement between parties: If party 1 sends a remote request structured a particular way, this is how party 2's software will respond.

API examples: - Google Maps API: allows developers embed geo-location data using JavaScript. The Google Maps API is designed to work on mobile and desktop. - YouTube API: allows developers integrate YouTube videos and functionalities into websites or applications. - Google Analytics API: allows to track website performance in terms of audience, monetization and other important metrics through the Google Analytics interface. The website the thesis come from has this implementation working.

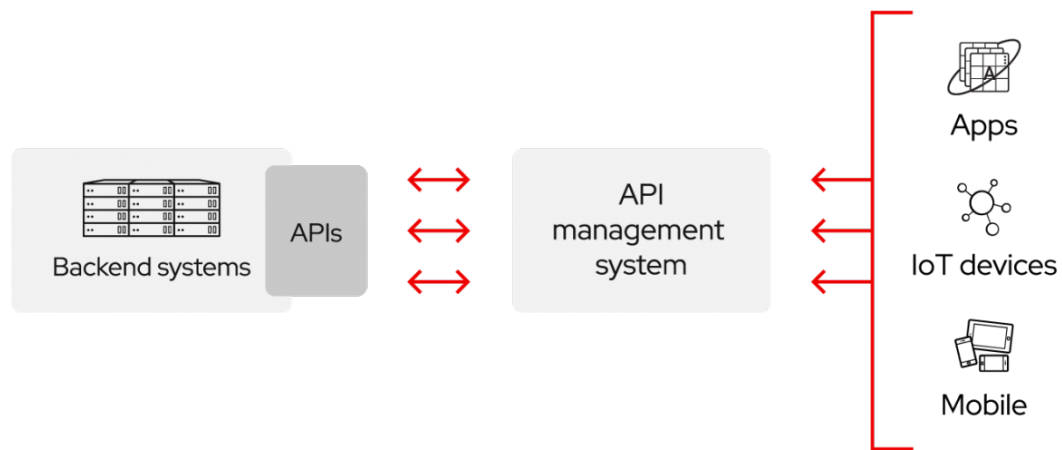


Figure 3.6: API functioning

Because APIs simplify how developers integrate new application components into an existing architecture, they help business and IT teams collaborate. Business needs often change quickly in response to ever shifting digital markets, where new competitors can change a whole industry with a new app. In order to stay competitive, it's important to support the rapid development and deployment of innovative services. Cloud-native application development is an identifiable way to increase development speed, and it relies on connecting a microservices application architecture through APIs.

3.4.1 What in practice an API does

API, strictly talking to this thesis extent, are cloud infrastructure that given a specific URL (and the most of the times credentials) are outputting data. Data comes in a given format (the most of the times JSON),so that they can

be instantly pre-processed and elaborated, in this case by one other software, the Shiny app.

(va prafrasato)

Since websites are continuously changed for many reasons API philosophy results in a smarter choice since it makes easy to access data and flex API endpoints to the business needs.

3.4.2 Plumber API

Plumber (api, 2020) allows the user to create a web API by simply adding decoration comments to the existing R code. Decorations are a special type of comments that suggests the plumber where and when the API parts are. Here below a toy example from the reference:

```
# plumber.R file

## Echo back the input
## @param msg The message to echo
## @get /echo
function(msg="") {
  list(msg = paste0("The message is: ", msg, ""))
}

## function(msg="") {
##   list(msg = paste0("The message is: ", msg, ""))
## }

## Plot a histogram
## @png
## @get /plot
function() {
```



```

  rand = rnorm(100)
  hist(rand)
}

```

```

## function() {
##   rand = rnorm(100)
##   hist(rand)
## }

```

```

## Return the sum of two numbers
## @param a The first number to add
## @param b The second number to add
## @post /sum
function(a, b) {
  as.numeric(a) + as.numeric(b)
}

```

```

## function(a, b) {
##   as.numeric(a) + as.numeric(b)
## }

```

this chunk of code assembled into a .R file has three endpoints

Special comments are marked as this `##` and they are followed by specific keywords denoted with `@`. - the name of the API (unique without the `@`) - the `@params` keyword refers to parameter that has to be inputted to give the result of the function define below. If in the function below default parameters are stated then the API response is the elaboration of the functions with those parameters. If the function does not specify any parameter, see the below endpoint below, plumber has to make sure that the function can run without them. - the `## @png` chunk piece specify the extension of the output file. - the `## @get /plot` decorations specify the type of HTTP request we are sending,

in this case a GET request. The user in this case is requesting some information to the API, the name is plot, so the expectationa are that a plot is given as response

Chapter 4

Methodologies

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Praesent sollicitudin arcu eget nulla convallis, in sagittis metus tincidunt. Integer vitae erat convallis, lobortis nunc id, commodo ipsum. Nam sem sem, tristique vitae dolor eget, laoreet aliquet libero. Nulla non feugiat diam. Ut vehicula, ante vitae rhoncus volutpat, eros orci viverra sem, ac tincidunt sem ex at sapien. Nam et odio condimentum, viverra nisl vitae, tempor nisl. Integer euismod convallis augue quis iaculis. Nulla at leo non nisi sollicitudin molestie vel at purus.

Nam tincidunt tellus et mattis luctus. Vivamus quis velit at nunc fermentum cursus. Duis elementum in nibh quis luctus. Suspendisse eget sem sit amet quam mattis egestas. Morbi ullamcorper metus eu dolor dapibus, id ultrices tellus euismod. Suspendisse malesuada felis vel tincidunt ullamcorper. Proin placerat auctor urna vel finibus. Sed non dictum orci. Ut volutpat pretium massa, in iaculis mi consectetur quis. Curabitur vitae condimentum nisi, sit amet consectetur justo. Curabitur non fermentum diam. Pellentesque vehicula laoreet elementum. Donec non porttitor ante, ut fermentum ante. In mollis consequat nisl euismod lobortis.

Vestibulum scelerisque dui eget est cursus, ut vestibulum libero pretium. Donec non viverra ligula. Suspendisse a viverra purus, in laoreet ligula. Nunc posuere libero ipsum, non vehicula massa gravida id. Interdum et malesuada fames ac ante ipsum primis in faucibus. Maecenas tortor risus, accumsan

vitae luctus a, molestie non augue. Nam felis arcu, volutpat vitae eros nec, sollicitudin aliquam orci. Sed sit amet consectetur dui. Morbi finibus, est at blandit consectetur, diam massa ultrices libero, ac dapibus nisl ex id urna. Donec venenatis, nibh malesuada ultricies varius, justo tellus maximus ante, vehicula viverra dolor lorem vitae lacus. Donec nec ultricies ex. Proin vehicula interdum ex a tincidunt.

Duis varius ornare velit, non finibus purus lacinia eget. Fusce fringilla arcu in mauris fermentum, at consequat mauris suscipit. Etiam cursus felis ut consequat maximus. Vestibulum auctor sollicitudin nisl, eget molestie enim faucibus sit amet. Aenean luctus non ligula scelerisque maximus. Aenean finibus, nulla sit amet interdum ornare, justo nisl tempor diam, sit amet sollicitudin lacus tortor sit amet mi. Aliquam erat volutpat. Quisque vehicula facilisis ligula ut porta. Aenean eleifend arcu luctus ligula congue lacinia. Proin nunc mauris, cursus vel risus at, dapibus imperdiet est. Fusce eget nisi nec eros vehicula ullamcorper. Integer feugiat vulputate lacus id posuere. Nunc tincidunt, ante ac convallis elementum, ipsum elit rutrum risus, vitae dapibus augue leo in lacus.

Integer et aliquam mi, ac blandit arcu. Integer sit amet tristique orci. Aenean consectetur tempor diam, id finibus massa tempus non. Cras sagittis nibh vitae aliquet laoreet. Sed facilisis luctus mi. Donec est sapien, porta consectetur varius in, pulvinar sed ligula. Suspendisse potenti. Maecenas porta neque at accumsan eleifend. Proin porta imperdiet ipsum, sed viverra purus. Interdum et malesuada fames ac ante ipsum primis in faucibus. Mauris dapibus libero ac ultrices commodo. Donec dictum lectus id urna volutpat, vitae eleifend neque gravida. Mauris aliquam augue ac eros sagittis interdum. Sed dapibus placerat bibendum.

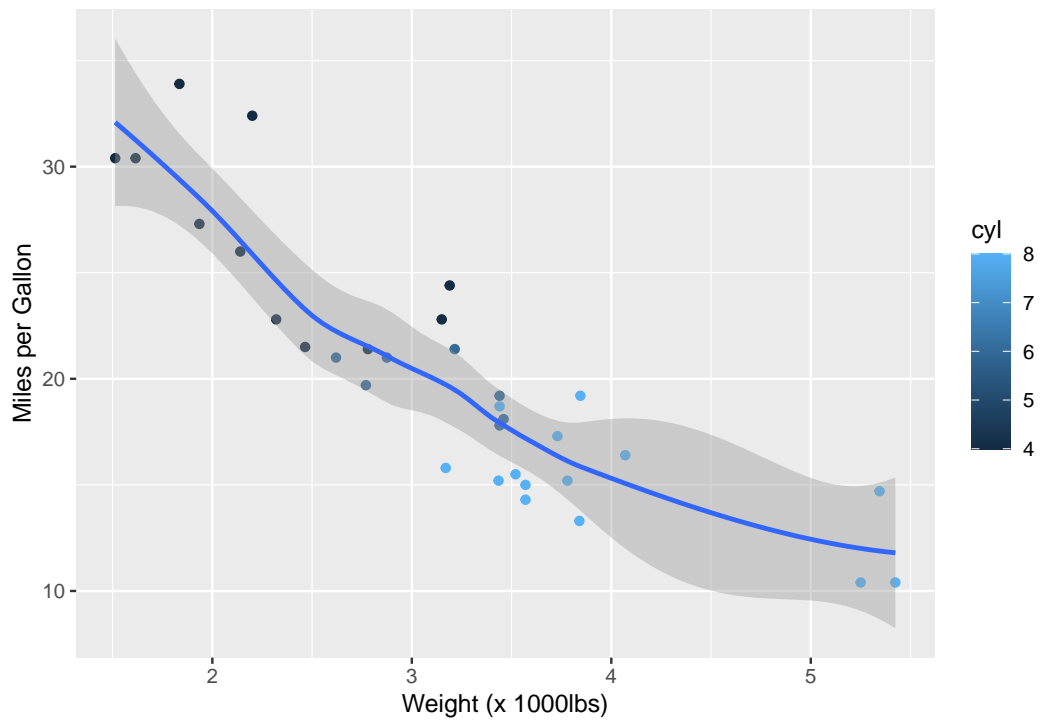
Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed imperdiet lacus ut maximus tempus. Integer dolor massa, ornare vitae ultrices pharetra, lacinia ut felis. Phasellus et vulputate mi, eu ultricies nisi. Phasellus interdum risus tristique tortor fermentum, blandit vehicula turpis elementum. Nullam eget

arcu faucibus, blandit lacus ut, pharetra felis. Praesent malesuada mattis augue vel convallis. Curabitur sodales sollicitudin nunc, eget gravida risus condimentum non. Cras laoreet eget erat eget dapibus. Vivamus volutpat consequat libero at dapibus. Nulla convallis sapien neque, sit amet scelerisque enim luctus at. Praesent vehicula convallis purus. In a metus nec tellus mattis vestibulum. Vestibulum ultrices eleifend quam eu tincidunt.

Phasellus mauris elit, volutpat sit amet suscipit ut, scelerisque sed quam. Quisque maximus, justo non porta sollicitudin, ante lorem placerat nulla, nec commodo mi mauris eget risus. Interdum et malesuada fames ac ante ipsum primis in faucibus. Vestibulum ornare eleifend sem quis mollis. Proin vestibulum tristique erat, eu suscipit neque ornare et. Aenean scelerisque, nibh nec rutrum varius, justo nulla ultrices erat, non maximus mi massa vel mi. Integer ante arcu, accumsan vitae quam nec, faucibus cursus leo. Donec feugiat nunc quis orci ornare dignissim. Vivamus posuere enim et odio pretium, non tincidunt diam finibus. Praesent ut odio vitae magna euismod tempus. Ut laoreet nibh quis iaculis faucibus. Vivamus neque turpis, tempus ac orci id, volutpat laoreet nunc.

```
library(tidyverse)
my_scatplot <- ggplot(mtcars, aes(x=wt, y=mpg, col=cyl)) +
  geom_point() +
  xlab('Weight (x 1000lbs)') +
  ylab('Miles per Gallon') +
  geom_smooth()
my_scatplot

## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



Nulla congue nunc elit, sit amet convallis urna lacinia eget. In massa erat, tempor eu erat ac, luctus sollicitudin felis. Curabitur urna ipsum, condimentum vel hendrerit a, imperdiet maximus ipsum. Morbi maximus malesuada efficitur. Phasellus id velit nec quam efficitur efficitur. Aenean a diam ut enim vehicula auctor ac nec sapien. Mauris in dapibus ex, quis consectetur ipsum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Curabitur at iaculis nulla. Aliquam turpis lacus, fermentum a malesuada eu, ornare sed libero. Fusce dictum fermentum ipsum, quis congue libero lacinia vel. Sed in finibus ligula. Vestibulum sagittis vel ex a mattis. Phasellus suscipit justo at augue porttitor pellentesque.

Nulla sed tellus vel nisl rutrum mattis a et est. Aenean dolor nunc, placerat quis enim a, fermentum malesuada enim. Aliquam erat volutpat. In ut est non velit vehicula hendrerit. Pellentesque eu erat finibus, viverra velit ac, tempus odio. Suspendisse potenti. Curabitur vulputate metus non ligula maximus sollicitudin. Donec facilisis augue sed urna lobortis, eu tempus mi ultrices. Aliquam fringilla sagittis felis, ac commodo arcu porttitor at. Ut lorem ex, gravida a porttitor in, feugiat vel ex. Aliquam consequat finibus ante, sed

commodo sem. Ut ac mollis dui. Curabitur nec eros congue, pulvinar quam at, tincidunt dolor. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Phasellus eu elit vulputate, scelerisque felis ut, accumsan nunc.

In non odio et risus rhoncus viverra nec in odio. Ut viverra metus eget nisi molestie sodales. Vestibulum ligula felis, malesuada quis vehicula eu, bibendum finibus diam. Donec ut rhoncus odio, tristique sodales orci. Praesent a ex lectus. Maecenas tristique sapien lorem, vel lobortis lorem porta sodales. Suspendisse ligula justo, iaculis id lorem eu, consequat vulputate lectus. Nunc molestie lacus at elit tempor, at rhoncus justo hendrerit. Nullam luctus lectus ac orci interdum, vitae varius libero pellentesque. Nunc lacinia erat lectus. Sed sit amet venenatis quam. Mauris interdum mauris sem, sit amet interdum eros tincidunt a.

Nam at dolor dui. Praesent efficitur leo erat, id blandit neque ultrices non. Morbi eget dignissim eros. Praesent rhoncus maximus accumsan. Quisque et consectetur odio, id vehicula leo. Integer tempor diam augue, nec rhoncus ex suscipit id. Sed nec tempor tellus. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi feugiat tincidunt tortor quis efficitur. Maecenas pellentesque dapibus nisi, sed commodo augue. Fusce condimentum dignissim quam id feugiat. Mauris maximus ex vel enim viverra, eget placerat massa consectetur. Integer pellentesque finibus ipsum, quis vestibulum elit ultrices a. Vivamus nunc velit, lobortis at hendrerit non, laoreet nec urna.

Vestibulum ullamcorper, lacus sed malesuada porta, lectus nisi lacinia augue, at mollis lorem purus sit amet metus. Quisque et mauris leo. Donec id risus id nisl auctor gravida. Suspendisse sed tempor risus. Integer ornare sem quis turpis accumsan finibus. Morbi in ex dui. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nam id laoreet erat. Curabitur tempor faucibus turpis, a bibendum turpis dignissim tempor.

Chapter 5

Applications

Some *significant* applications are demonstrated in this chapter.

5.1 Example one

5.2 Example two

Chapter 6

Final Words

We have finished a nice book.

Bibliography

- (2020). 7.2. advantages of using docker red hat enterprise linux 7.
- (2020). An api generator for r.
- (2020). Get docker.
- (2020). User agent: Learn your web browsers user agent now.
- Allaire, J., Xie, Y., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., Wickham, H., Cheng, J., Chang, W., and Iannone, R. (2020). *rmarkdown: Dynamic Documents for R*. R package version 2.3.
- Chang, W., Cheng, J., Allaire, J., Xie, Y., and McPherson, J. (2020). *shiny: Web Application Framework for R*. R package version 1.4.0.2.
- Inc, F., Weststrate, M., Russell, K., and Dipert, A. (2020). *reactR: React Helpers*. R package version 0.4.3.
- Lovelace, R., Nowosad, J., and Muenchow, J. (2019). *Geocomputation with R*. CRC Press.
- Trestle Technology, LLC (2018). *plumber: An API Generator for R*. R package version 0.4.6.
- was previously an Economist at the Indeed Hiring Lab, A. F. F. (2020). Indeed tech skills explorer: Today’s top tech skills.
- Wikipedia (2020). Web scraping — wikipedia, l’enciclopedia libera. [Online; in data 15-luglio-2020].

Wikipedia contributors (2020). Cron — Wikipedia, the free encyclopedia. [Online; accessed 15-September-2020].

Wikiversità (2020). Scheduling — wikiversità,. [Online; accesso il 15-settembre-2020].

Xie, Y. (2016). *bookdown: Authoring Books and Technical Documents with R Markdown*. Chapman and Hall/CRC, Boca Raton, Florida. ISBN 978-1138700109.

Xie, Y. (2020). *bookdown: Authoring Books and Technical Documents with R Markdown*. R package version 0.20.