

UNIVERSITÀ CATTOLICA SACRO CUORE

STATISTICAL AND ACTUARIAL SCIENCES

MJ: DATA BUSINESS ANALYTICS

REST API for Real Estate rental data, a spatial Bayesian modeling approach with INLA

Author:
Niccolò SALVINI

Supervisor:
Dr. Marco DELLAVEDOVA

Assistant Supervisor:
Dr. Vincenzo NARDELLI

AY 2019 / 2020



REST API for Real Estate rental data, a spatial Bayesian modeling approach with INLA.

Niccolò Salvini¹

Last compiled on 05 novembre, 2020

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

Contents

1	Introduction	7
2	Scraping	9
2.1	What is Web Scraping	10
2.1.1	Immobiliare.it Webscraping website structure	14
2.1.2	Immobiliare.it Webscraping content architecture with rvest	15
2.2	Scraping Best Practices and Security provisions	24
2.3	Security provisions: User Agents, Proxies and Handlers	26
2.3.1	User Agents Spoofing	28
2.3.2	Handlers and Trycatches	29
2.4	Parallel Computing	31
2.5	Open Challenges	35
2.6	Legal Profiles (ancora non validato)	36
3	REST API Infrastructure	38
3.1	Scheduler	40
3.1.1	Cron Jobs	41
3.2	Docker	43
3.2.1	Why Docker	44
3.2.2	Dockerfile	46
3.3	REST API	48
3.3.1	Plumber REST API	50
3.3.2	Immobiliare.it REST API	53
3.3.3	REST API documentation	55
3.4	NGINX reverse proxy server	56
3.5	AWS EC2 server	56

3.5.1	Launch an EC2 instance	57
3.6	Further Integrations	58
4	INLA computation	59
4.1	Latent Gaussian Models LGM	60
4.2	Approximation in INLA setting	64
4.2.1	further approximations (prolly do not note include) . .	65
4.3	R-INLA package in a bayesian hierarchical regression perspective	67
4.3.1	Overview	67
4.3.2	Linear Predictor	68
5	Point Referenced Data Modeling	72
5.1	Gaussian Process (GP)	75
5.2	Spatial Covariance Function	78
5.2.1	Matérn Covariance Function	80
5.3	Hedonic models Literature Review and Spatial Hedonic Price Models	82
5.4	Point Referenced Regression for univariate spatial data	85
5.5	Hierarchical Bayesian models	87
5.6	INLA model through spatial hierarchical regression	90
5.7	Spatial Kriging	91
5.8	Model Checking	92
5.9	Prior Specification	94
6	SPDE approach	95
6.1	Set SPDE Problem	96
6.2	SPDE within R-INLA	97
6.3	First Point Krainsky Rubio TOO TECHNICAL	97
7	Exploratory Analysis	99
7.1	Data preparation	102
7.1.1	NA Removal and Imputation	103
7.2	Spatial Autocorrelation assesement	105
7.3	Model Specification	106
7.4	Mesh building	106

<i>CONTENTS</i>	3
7.4.1 Shinyapp for mesh assessment	107
7.4.2 BUilding SPDE model on mesh	108
7.5 Spatial Kriging (Prediction)	108
8 Model Selection & Fitting	109
8.1 Model Criticism	109
8.2 Spatial Kriging	109
8.3 Model Checking	109
9 Final Words	110

List of Tables

List of Figures

2.1	general website structure	11
2.2	html_tree	13
2.3	immobiliare.it website structure, author source	15
2.4	rvest flow chart, missing source	16
2.5	immobiliare.it content structure, author's source	17
2.6	How Web Works	27
2.7	computational complexity analysis with Furr	32
2.8	computational complexity analysis with Furr	34
3.1	complete infrastructure (Matt Dancho source)	40
3.2	crontab	42
3.3	docker container vs VM	43
3.4	docker-stats	45
3.5	dockerfile	47
3.6	API functioning	51
3.7	swagger	54
3.8	aws_dashboard	57
4.1	CCD to spdetoy dataset, source Marta Blangiardo (2015) . . .	66
4.2	summary table list object, source: Krainski (2019)	68
4.3	SPDEtoy plot, author's source	69
4.4	linear predictor marginals, author's creation	71
5.1	point referenced data example, Milan Rental Real Estate, Au- thor's Source	73
5.2	3D scatterplot and surface, Stockton data.	75
5.3	variogram example	77

5.4	isotropy VS anisotropy, source Blanchet-Scalliet et al. (2019) .	78
5.5	matern correlation function for 4 diff values of nu with phi fixed, author's source	81
5.6	9 levels cat vs observaitons, source Marta Blangiardo (2015) .	88
5.7	Spatial prediction representation through DAG, source Marta Blangiardo (2015)	92
6.1	lattice 2D regular grid	98
7.1	Missingness Heatmap plot	104
7.2	Missingness co-occurrence plot	105
7.3	Traingularization intuition, Krainski (2019) source	106

Chapter 1

Introduction

Trento:

- Argomento
- Problema
- Obiettivi
- Metodo
- Struttura della tesi

Main themes:

- Research Question
- Milan Real Estate Controversies in relation to research question
- why the API (perchè mi mancano i dati e perchè è il futuro)
- Open Data discussion personal hope of data sharing and benefits from open source
- Why a Bayesian approach
- Why INLA

As a general discussion technologies implied can be thought as the distance between a service running locally on a laptop and something that it can actually be put into production, shared among company stakeholders, solving

business related problems. When such technologies are applied data scientist and interlocutors gradually close the gap. Insights are better communicated, data is up-to-date and automation can save time. Nonetheless when the infrastructure is structured with vision then integrating or substituting existing technologies is not trivial. Anyway technologies can not be always embedded because they might be exclusively designed to work only on certain back ends, therefore some choices are not into discussion. With foresight RStudio by setting future-oriented guidelines has spent a lot of effort giving its users an easy, integrated and interconnected environment. By that it is meant that the RStudio community has tried to either integrate or open the possibility to a number of technologies that fill the blanks in their weaker parts. On top of many, an entire package has been dedicated to democratize REST APIs (Plumber (Trestle Technology, LLC, 2018)). As a further example developers in RStudio have created an entire new paradigm i.e. Shiny (Chang et al., 2020), a popular web app development package, that enforces the developer to have front-end and back-end technologies tied up in the same IDE. They also added performance monitoring and optimization packages that are fitted into shiny such as shinytest [metti tag] and shinyloadtest [metti tag] to simulate sessions and verify network traffic congestion.

Chapter 2

Scraping

The following chapter covers a gentle introduction and concepts of web scraping centered on immobiliare.it. It starts in general terms by ideally segmenting websites in two meta-concepts: website structure and content architecture, then it provokes this segmentation to the specific immobiliare.it case. The abstract workflow will ease the identification of the first “high level” challenge to be fronting during scraping, which is mainly understanding how the website is structured and how to reverse engineer the url composition. The structure is unrolled so that each part can be singularly and directly accessed and the problem can be passed to “lower levels”. At this point a rooted-tree graph representation is used to map scraping functions into immobiliare content architecture. By means of **rvest** scraping is possible and the main function involved are presented. A specific function to scrape price (the response var of the analysis) is shown and then also a second function, that takes care of grouping all the scraping functions into a single one. Scraping best practices are applied both on the web server side, kindly requesting permission and delayed sending rate; and from the web client side by granting continuous scraping, avoiding server blocks by User Agent rotation and trycatches / handlers for easy debugging. Then run time issues are fronted presenting two different options for looping construction, **furrr** and **foreach**. The latter has displayed interesting results and further improvements are taken into

consideration. Then an overview of the open challenges is offered so that this work might be extended or integrated with other technologies. In the end legal profiles are addressed comparing scraping results and difficulties with a counterpart case study.

2.1 What is Web Scraping

Web Scraping is a technique aimed at extracting data from static or dynamic internet web pages. It can be applied simultaneously or automatically by a scheduler that plans execution at a given time. Content in webpages is the most of the times well organized and accessible. (puoi dire meglio) This is made possible by the effort put into building both the *website structure* and the *content architecture*. For website structure it is meant the way urls, pointing to webpages, are arranged throughout the website. Website structure constitutes a *first dimension* of hierarchy. Some popular structure examples might regard social-networks where posts can be scrolled down within a single url page named the wall. Scrolling down might end due to few posts, but the perception related to the scroll option is a never-ending webpage associated to a single url. Instead personal profiles are dedicated to a specified unique url and even in profiles posts are allocated into a sub domain and might be scrolled down arranged by time since the day of social subscription. Online newspapers display their articles in the front wall and by accessing to one of them all the related following articles sometimes can be reached by an arrow, pointing right or left. Articles can also be suggested and the more the website is explored the more articles are likely to be seen twice within the same session. It will soon end up into a suggestion loop that it will recursively shows the same contents; recursive structures are popular in newspaper-type websites. Online Retailers as Amazon, based on filters, groups inside a single webpage (i.e. page n° 1) a fixed set of items, having their dedicated personal url attached to them. Furthermore Amazon offers the opportunity to skip to the following page (i.e. page n° 2), searching for another different and fixed set of items

and so on until the last page/url. Generally website structures try to reflect both the user expectations with respect to the product on the website and the design expression of the web developer. For these reasons for each websites usage category exists a multitude of content architectures. For each content architectures there are multiple front end languages which are destined to multiple end users. In the future expectations are tailor made webpages on users with respect to its personal preferences. Moreover web design in scraping plays an important role since the more sophisticated graphical technologies are implied, the harder will be to scrape information.

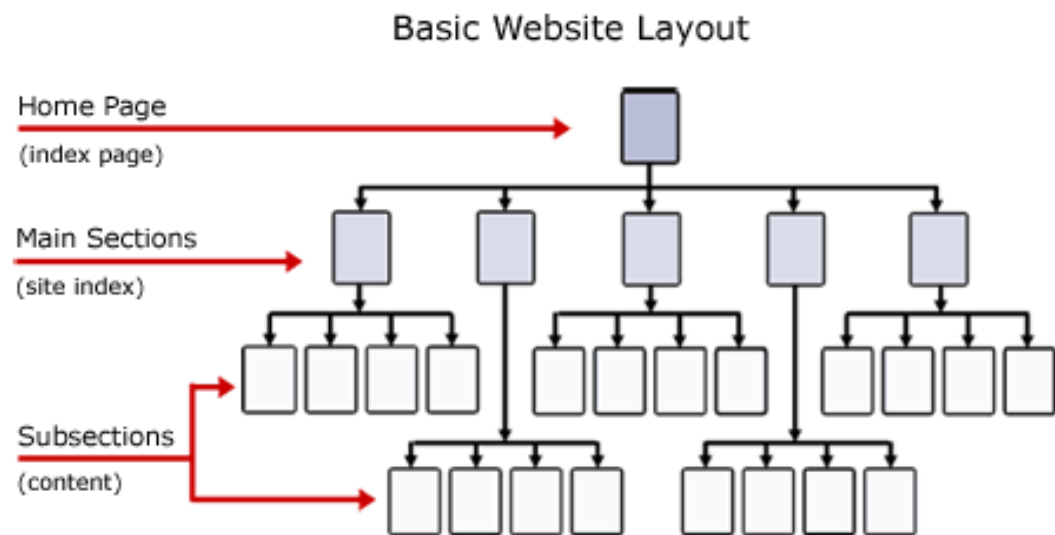


Figure 2.1: general website structure

A *second dimension* of hierarchy is brought by content architecture in the name of the language used for content creation and organization i.e. HTML. HTML stands for Hyper Text Markup Language and ... HTML drives the hierarchy structure that is then generalized to the website structure. According to this point of view the hierarchical website structure is a consequence of the content architecture by means of HTML language (*arborescence*: direction from root to leaves). CSS language stands for Cascading Style Sheets and takes care of the style of the webpage. The combination of HTML and CSS offers a wide flexibility in building web sites, once again expressed by the vast amount of different designs on the web. Some websites' components also might be tuned

by Javascript language, which in the context of scraping adds a further layer of difficulty. As a matter of fact since Javascript components are dynamic within the webpage, scraping requires specialized libraries to enable different parser to get the content. CSS allows the scraper to target a class of objects in the web page that shares same style (e.g. same css query) so that each element that belongs to the class (i.e. share same style) can be gathered. This practice provides enormous advantages since by CSS a set of objects can be obtained within a single function call. First and Second dimension of the scraping problem imply hierarchy, a simple way to approach the problem is to represent it through already known data structures. One way to imagine hierarchy in both of the two dimensions are graph based data structures named as **Rooted Trees**. By analyzing the first dimension through the lenses of Rooted trees it is possible to compress the whole problem into the general graph based jargon (Diestel, 2006). Rooted trees must start with a root node which is the domain of the web page. Each *Node* is a url destination and each *Edge* is the connection between nodes. Connections have been made possible in the website by nesting urls inside webpages so that within a single webpage the user can access to a number of other related links. Furthermore, as an extension to the rooted tree framework a general graph theory component is introduced, i.e. the *Weight*. Each edge is associated to a weight whose interpretation is the run time cost to walk from a node to its connected other nodes (e.g. from a url to the other). In addition the content inside each node takes the name of payload, which ultimately is the scope of the scraping processes. The walk from node 17 to node 8 in figure below (even though that is the case of a binary rooted tree) is called path and it represented as an ordered list of nodes connected by edges. In this context each node can have both a fixed and variable outgoing sub-nodes that are called *Children*. When root trees have a fixed set of children are called *k-ary* rooted trees. A node is said to be *Parent* to other nodes when it is connected to them by outgoing edge, in right figure below “head” is the parent of nodes “title” and “meta”. Nodes in the tree that shares the same parent node are said *Siblings*, “head” and “body” are siblings in figure

@ref(tree_html). Moreover *Subtrees* are a set of nodes and edges comprised of a parent and its descendants e.g. node “body” with all of its descendants might constitute a subtree. The concept of subtree in both of the dimensions plays crucial role in cutting run time scraping processes and fake headers provision (see section 2.3.1). If the website structure is locally reproducible and the content architecture within webpages tends to be equal, then functions for a single subtree might be extended to the rest of others subtrees that are siblings to the same parent node. Local reproducibility is a property according to which starting from a single url all the related urls can be inferred from a pattern. Equal content architecture throughout different single links means to have sort of standard shared-within advs criteria that each single rental advertisement has to refer. In addition two more metrics have might better describe the tree: *level* and *height*. The level of a node **L** counts the number of edges on the path from the root node to **L**. The height is the maximum level for any node in the tree, from now on **H**. What is worth to be anticipating is that functions are not going to be applied directly to siblings in the more general rooted tree. Instead it would be better segmenting the highest level rooted tree into a sequence of single subtrees whose roots are the siblings for reasons explained in section 2.3.1.

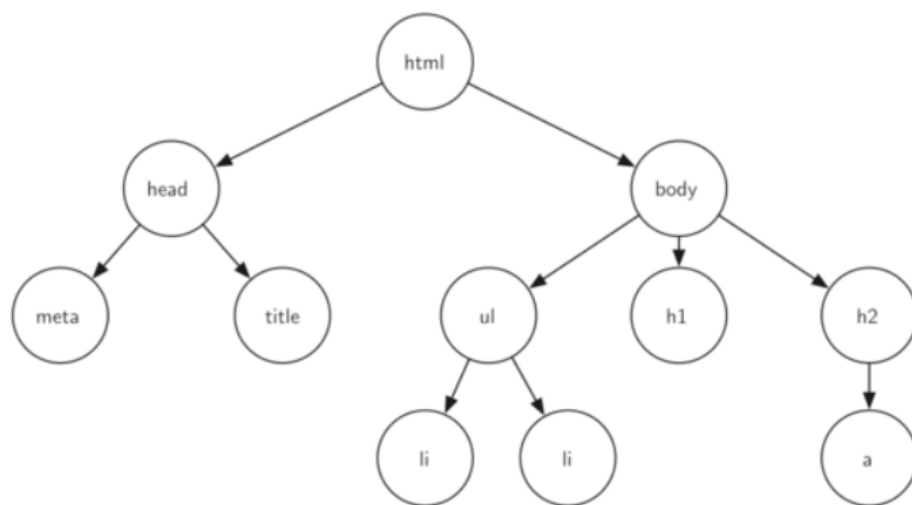


Figure 2.2: html_tree

2.1.1 Immobiliare.it Webscraping website structure

The structure of the website resembles the one from the popular online retailer Amazon. According to the query filters selected in a dedicated section (e.g. city, number of rooms 5, square footage less than 60 m^2 etc), the url is shaped so that each further filter is appended at the end of the domain url `https://www.immobiliare.it/` root node. Filters are appended to the domain with a proper syntax, not all the syntax are equal. Once filters are all applied to the root domain this constitutes the new url root domain node that might have this appearance by filtering for rents in Milan city when square footage is less than 60 m^2 : `domain+affitto-case/milano/?superficieMinima=60`. Since this is true only for page n°1 containing the first 25 advs (see figure @ref(fig:website_tree1)) all the remaining siblings nodes corresponding to the subsequent pages have to be generated. Here the utility of Local reproducibility property introduced in the previous section. The remaining siblings, i.e. the ones belonging to page 2 (with the attached 25 links), to page 3 etc. can be generated by adding a further filter `&pag=n`, where n is the page number reference (from now on referred as *pagination*). Author customary choice is to stop pagination up to 300 pages since spatial data can not be too large due to computational constraints. Up to this process pagination has generated a list of siblings nodes whose children elements number is fixed (i.e. 25 links per page @ref(fig:website_tree1) lower part). That makes those trees *k-ary*, where k is 25 indicating the number of children leaves. K-ary trees are rooted trees in which each node has no more than k children, in this particular case final leaves. The well known binary rooted tree is actually a special case of k-ary when $k = 2$. Filters reverse engineering process and 25-ary trees with equal content structure across siblings allow to design a single function to call that could be mapped for all the other siblings. In addition in order to further unroll the website a specific scraping function grabs the whole set of 25 links per page. As a result a single function call of `scrape_href()` can

grab the links corresponding to page 1. Then the function is mapped for all the generated siblings nodes (i.e. up to 300) obtaining a collection of all links belonging to the set of pages. Ultimately the complete set of links corresponds to every single advertisement posted on immobiliare.it at a given time.

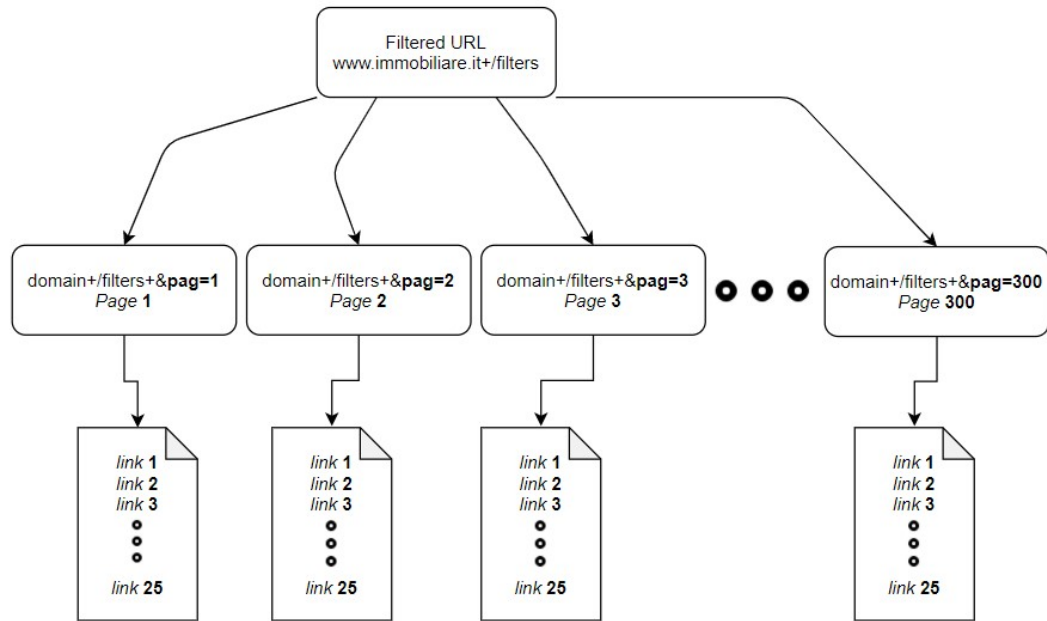


Figure 2.3: immobiliare.it website structure, author source

2.1.2 Immobiliare.it Webscraping content architecture with `rvest`

To start a general scraping function the only requirement is a target url (i.e. the filtered root node url). Then a list `html_session` object is opened by specifying the url and the request data that the user need to send to the web server (see left part to dashed line image 2.4). Information to be attached to the web server request will be further explored later, though they are mainly three: User Agents, emails references and proxy servers. `html_session` objects contains a number of useful information such as: the url, the response, cookies, session times etc. Once the connection is established (request response 200) all the following operations rely on the opened session, in other words for the

time being in the session the user will be authorized with the before-provided characteristics through the request. The list object contains mostly the html content of the webpage and that is where data needs to be parsed. The list can disclose as well other interesting meta information related to the session but they are not collected. The figure in @ref(fig:rvest flow chart, missing source) sketches what the scraping needs to front:

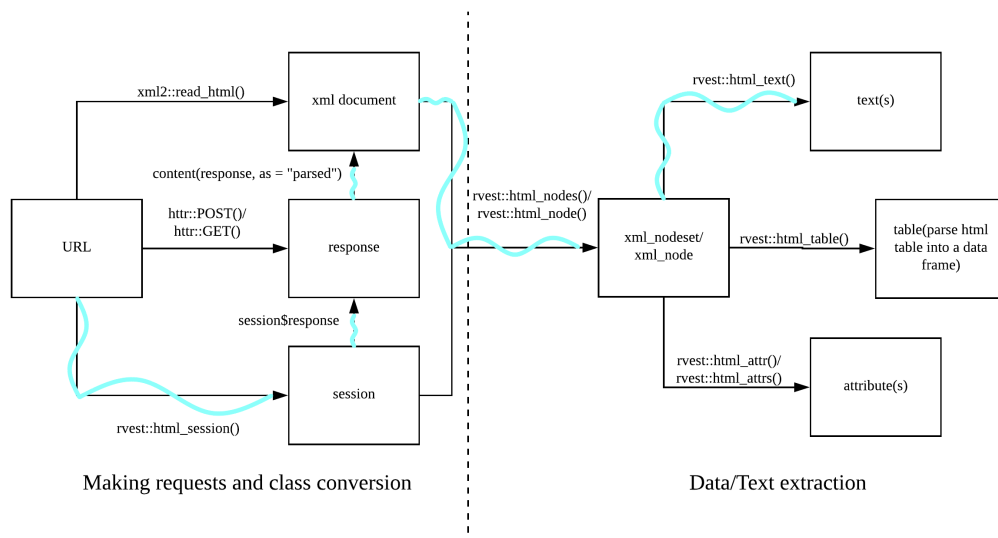


Figure 2.4: rvest flow chart, missing source

To the right of dashed line in the flow chart are painted a sequence of **rvest**(Wickham, 2019) functions that follow a general step by step text comprehension rules. **rvest** first handles parsing the html content of the web page within the session object **read_html()**, secondly It looks for a single node **html_nodes()** through a CSS query. CSS is a way to route **rvest** to consider a precise node in the webpage. Thirdly it converts the content (i.e. payload) into a human readable text with **html_text()**. The sketch of immoliare.it content structures is reported in the figure below:

The code chunk below shows a function that can scrape the price.

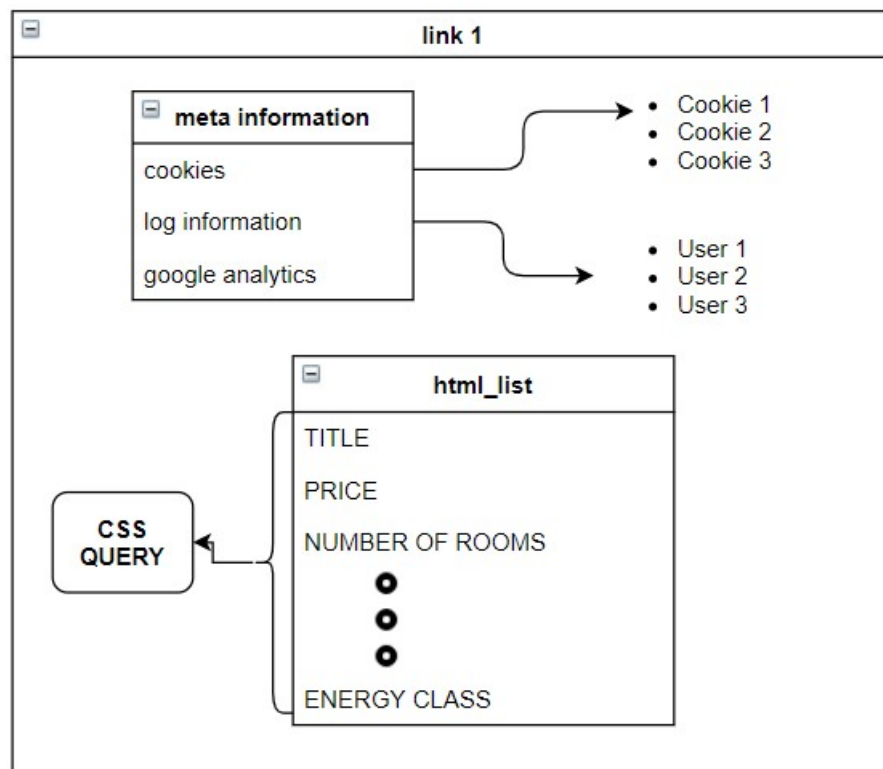


Figure 2.5: immobiliare.it content structure, author's source

```

scrapeprice.imm = function(session) {

  opensess = read_html(session)
  price = opensess %>% html_nodes(css = ".im-mainFeatures__title") %>% html_text() %>% str_trim()

  if (is.null(price) || identical(price, character(0))) {
    price2 = opensess %>% html_nodes(css = ".im-features__value , .im-features__value") %>% html_text() %>% str_trim()

    if ("prezzo" %in% price2) {
      pos = match("prezzo", price2)
      return(price2[pos + 1]) %>% str_replace_all(c(`\200` = "", `\\.` = "")) %>% str_extract("\\\\-*\\d+\\.\\d*") %>% str_replace_na() %>% str_trim()
    }
  }
}

```

```

        "Prezzo Su Richiesta")
    } else {
        return(NA_character_)
    }
} else {
    return(price) %>% str_replace_all(c(`\200` = "", `\\.` = "")) %>% str_
        str_replace_na() %>% str_replace("NA", "Prezzo Su Richiesta")

}

}

```

The function takes as a single argument a session object which is at first initialized in one other function. Then It reads the inner html content in the session storing the information into an obj called the **opensess**. Another obj is created, namely **price**, right after the pipe operator a css query into the html is called. The css query **.im-mainFeatures__title** points to a precise data location inside the web page, right below the main title. Expectation are that **price** is a one-element chr vector, containing the price and other unnecessary non-UTF characters. Then the algorithm enters into the first **if** statement. The handler checks if the object **price** is empty. If it doesn't the algorithm jumps to the end of the algorithm and returns the cleaned quantity. But If it does it takes again the **opensess** and tries with a second css query **.im-features__value** , **.im-features__title** where price might be also found. Please note that this is all done within the same session, so no more additional request information has to be sent. Since the latter css query points to data stored inside a list, for the time being the newly created obj **price2** is a list containing various information. Then the algorithm flow enters into the second **if** statement that checks whether the **"prezzo"** is matched the list or not, if it does it returns the +1 position index element with respect to the "prezzo" positioning. This happens because data in **price2** list are stored

by couples sequentially, e.g. [title, “Appartamento Sempione”, energy class, “G”, “prezzo”, 1200/al mese]. When it returns the element corresponding to +1 position index it applies also some data wrangling with `stringr` package to keep out overabundant characters. The function then escapes in the else statement by setting `price2 = NA_Character_` once no css query could be finding the price information. the *character-string* type has to be imposed due to fact that later they can not be bind. In other words if the function is evaluated for a url and returns the price quantity, but then is evaluated for url2 and outputs NA (no character) then results can not be combined into dataframe due to different types.

Once all the functions have been created they need to be called together and then data coming after them need to be combined. This is done by `get,data.catsing()` which at first checks the validity of the url, then takes the same url as input and filters it as a session object. Then simultaneously all the functions are called and then combined. All this happens inside a `foreach` parallel loop called by `scrape.all.info()`

```
scrape.all.info = function(url = "https://www.immobiliare.it/affitto-case...",
                           vedi = FALSE,
                           scrivi = FALSE,
                           silent = FALSE){

  if (silent) {
    start = as_hms(Sys.time()); cat('Starting the process...\n\n')
    message('\nThe process has started in',format(start,usetz = TRUE))
  }

  # open parallel multisession
  cl = makeCluster(detectCores()-1) #using max cores - 1 for parallel process
  registerDoParallel(cl)
  start = as_hms(Sys.time())

  if (silent) {
```

```

    message('\n\nStart all the requests at time:', format(start, usetz = T))
}

ALL = foreach(i = seq_along(links),
              .packages = lista.pacchetti,
              .combine = "bind_rows",
              .multicombine = FALSE,
              .export = "links" ,
              .verbose = TRUE,
              .errorhandling='pass') %dopar% {
  source("utils.R")
  sourceEntireFolder("functions_singolourl")
  get.data.catsing = function(singolourl){

    # dormi()
    #
    if(!is_url(singolourl)){
      stop(paste0("The following url does not seem either to e
    })

    session = html_session(singolourl, user_agent(agents[sampl
    if (class(session) == "session") {
      session = session$response
    }

    id      = tryCatch({scrapehouse.ID(session)},
                      error = function(e){ message("some p
    lat     = tryCatch({scrapelat.imm(session)},
                      error = function(e){ message("some p
    long    = tryCatch({scrapelong.imm(session)},
                      error = function(e){ message("some p
    location = tryCatch({take.address(session)},

```

```

                                error = function(e){ message("some p
condom      = tryCatch({scrapecondom.imm(session)},
                                error = function(e){ message("some p
buildage    = tryCatch({scrapeagebuild.imm(session)},
                                error = function(e){ message("some p

...

combine = tibble(ID          = id,
                  LAT         = lat,
                  LONG        = long,
                  LOCATION    = location,
                  CONDOM      = condom,
                  BUILDAGE    = buildage,

...

return(combine)
}

stopCluster(cl)
return(ALL)
}

```

The skeleton constitutes a standard format adopted for many other scraping function in the analysis. Being equal the css query what it changes is the matching term, i.e. “numero camere” instead of “prezzo” to look for how many rooms there are in the house. This is true for all the information contained in the list accessed by the fixed css query. Those that are not they are a few and they do not need to be scraped. In addition some other functions outputs need to undergo to further heavy cleaning steps in order to be usable As

a consequence of that functions need also to be broken down by pieces into many single .R files whose names correspond to each important information. Below it is printed the tree structure folder that composes the main elements of the scraping procedure. It can be noticed that the two folders, namely functions_singolourl and functions_url enclose all the single functions that allow to grab single information from session. Folders with a customized function are then sourced within the two main functions, scrape.all and scrape.all.info so data can be extracted.

```

levelName
1 immobiliare.it-WebScraping
2   |--functions_singolourl
3   |   |--0scrapesqfeetINS.R
4   |   |--0scrapenroomINS.R
5   |   |--0scrapepriceINS.R
6   |   |--0scrapetitleINS.R
7   |   |--ScrapeAdDate.R
8   |   |--ScrapeAge.R
9   |   |--ScrapeAgeBuilding.R
10  |   |--ScrapeAirConditioning.R
11  |   |--ScrapeAptChar.R
12  |   |--ScrapeCatastInfo.R
13  |   |--ScrapeCompart.R
14  |   |--ScrapeCondom.R
15  |   |--ScrapeContr.R
16  |   |--ScrapeDisp.R
17  |   |--ScrapeEnClass.R
18  |   |--ScrapeFloor.R
19  |   |--ScrapeHasMulti.R
20  |   |--ScrapeHeating.R
21  |   |--ScrapeHouseID.R

```



```
22 | |--ScrapeHouseTxtDes.R
23 | |--ScrapeLAT.R
24 | |--ScrapeLONG.R
25 | |--ScrapeLoweredPrice.R
26 | |--ScrapeMetrature.R
27 | |--ScrapePhotosNum.R
28 | |--ScrapePostAuto.R
29 | |--ScrapePropType.R
30 | |--ScrapeReaReview.R
31 | |--ScrapeStatus.R
32 | |--ScrapeTotPiani.R
33 | |--ScrapeType.R
34 | °--take_location.R
35 |--scrapeALL.R
36 |--scrapeALLINFO.R
37 |--functions_url
38 | |--ScrapeHREF.R
39 | |--ScrapePrice.R
40 | |--ScrapePrimaryKey.R
41 | |--ScrapeRooms.R
42 | |--ScrapeSpace.R
43 | °--ScrapeTitle.R
44 |--libs.R
45 |--utils.R
46 |--README.Rmd
47 |--README.md
48 °--simulations
49 | |--rt_match_vs_forloop.R
50 | °--runtime_simul.R
```

2.2 Scraping Best Practices and Security provisions

Robots.txt files are (rivedi citation) a way to kindly ask webbots, spiders, crawlers, wanderers and the like to access or not access certain parts of a webpage. The de facto ‘standard’ never made it beyond a informal “Network Working Group INTERNET DRAFT”. Nonetheless, the use of robots.txt files is widespread (e.g. <https://en.wikipedia.org/robots.txt>, <https://www.google.com/robots.txt>) and bots from Google, Yahoo and the like will adhere to the rules defined in robots.txt files, although their *interpretation* of those rules might differ.

Robots.txt files are plain text and always found at the root of a website’s domain. The syntax of the files in essence follows a fieldname: value scheme with optional preceding user-agent: ... lines to indicate the scope of the following rule block. Blocks are separated by blank lines and the omission of a user-agent field (which directly corresponds to the HTTP user-agent field) is seen as referring to all bots. # serves to comment lines and parts of lines. Everything after # until the end of line is regarded a comment. Possible field names are: user-agent, disallow, allow, crawl-delay, sitemap, and host. For further notions (Meissner and Ren, 2020, Google (2020))

Some interpretation problems:

- finding no robots.txt file at the server (e.g. HTTP status code 404) implies that everything is permitted
- subdomains should have there own robots.txt file if not it is assumed that everything is allowed
- redirects involving protocol changes - e.g. upgrading from http to https - are followed and considered no domain or subdomain change - so whatever is found at the end of the redirect is considered to be the - robots.txt file for the original domain

- redirects from subdomain www to the domain is considered no domain change - so whatever is found at the end of the redirect is considered to be the robots.txt file for the subdomain originally requested

For the thesis purposes it has been designed a dedicated function to assess whether the domain or the related paths require specific actions or they prevent some activity on the target. The following `checkpermission()` function has been integrated inside the scraping system and it is called once at the starting point.

```
dominio = "immobiliare.it"

checkpermission = function(dom) {

  robot = robotstxt(domain = dom)
  vd = robot$check()[1]
  if (vd) {
    cat("\nrobot.txt for", dom, "is okay with scraping!")
  } else {
    cat("\nrobot.txt does not like what you're doing")
    stop()
  }
}

## metti path allowed check
checkpermission(dominio)
```

```
##
## robot.txt for immobiliare.it is okay with scraping!
```

Further improvements in this direction might come with the `polite` package (Perepolkin, 2019) which combines the power of the `robotstxt`, the `ratelimitr` (Shah, 2018) to limit sequential requests together with the

`memoise` (Wickham et al., 2017) for response caching. This package is wrapped up around 3 simple but effective ideas:

The three pillars of a polite session are seeking permission, taking slowly and never asking twice.

The three pillars constitute the *Ethical* web scraping manifesto (Densmore, 2019) which are common shared practises that are aimed to self regularize scrapers. These have to be intended as best practices, not in any case as law enforcements, however many scrapers themselves, as website administrators or analyst, have fought in their daily working tasks with bots. Crawling bots in intensive scraping processes might fake real client navigation logs and as a consequence might induce distorted analytics. Due to this fact comes the need to find a common operating ground and therefore politely asking for permission.

2.3 Security provisions: User Agents, Proxies and Handlers

HTTP requests to the website server by web clients come with some mandatory information packed in it. The process according to which HTTP protocols allow to exchange information can be easily thought with an everyday real world analogy. As a generic person A rings the door's bell of person B's house. A comes to 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 requests. In certain cases it can even forbid the user log to the website. The language the two parties exchanges are coded in numbers that ranges from 100 to 511, each of which has its own specific 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 practices.

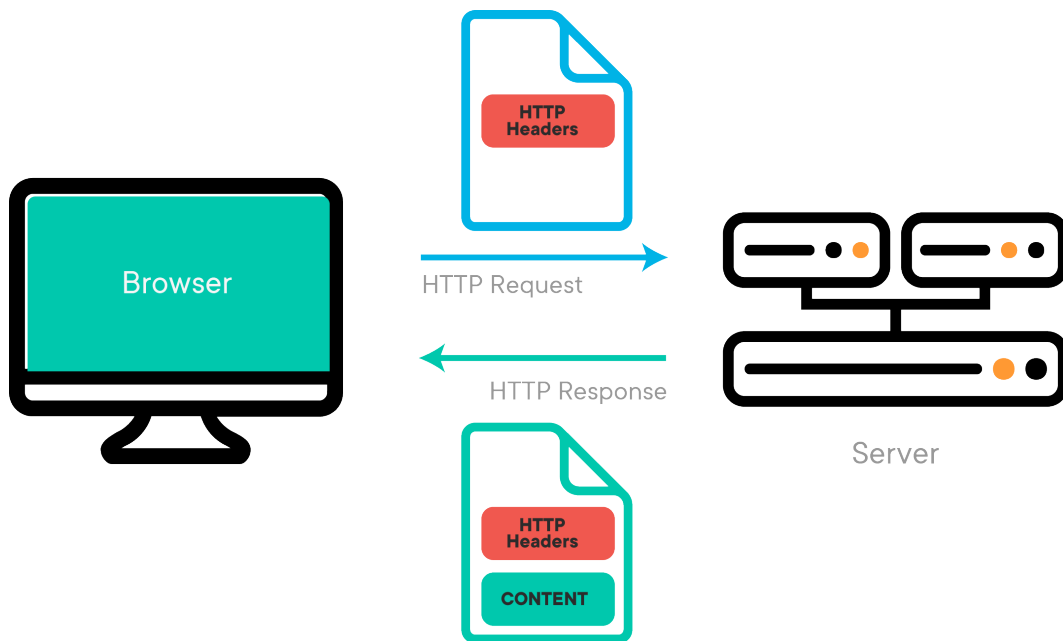


Figure 2.6: How Web Works

This procedure is a daily issue to people that are trying to collect information from websites. Google performs it everyday with its spider crawlers, which are very sophisticated bots that scrapes 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. Precautions have not been taken lightly, and a simple but effective approach is proposed.

2.3.1 User Agents Spoofing

A user agent (WhoIsHostingThis.com, 2020) is a string of characters in each web browser that serves as identification card. 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 includes the user application or software, the operating system (and their versions), the web client, the web client's version, as well as the web engine responsible for the content display (such as AppleWebKit). The user agent string is sent in the form of a HTTP request header. Since User Agents acts as middle man between the client request and the server response, then from a continuous scraping point of view it would be better rotating them, so that each time the middle man looks different. The solution adopted builds a vector of user agent strings identified by different specifications, different web client, different operating system and so on, then samples 1 of them. Then whenever a request from a web browser is sent to a web server, 1 random sample string is drawn from the user agents pool. So each time the user is sending the request it appears to be a different User Agent. Below the user agents rotation pool:

```
set.seed(27)
```

```
agents = c("Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like  

  \"Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko  

  \"Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like  

  \"Mozilla/5.0 (Macintosh; Intel Mac OS X 10_12_1) AppleWebKit/602.2.14 (KHTML  

  \"Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko  

  \"Mozilla/5.0 (Macintosh; Intel Mac OS X 10_12_1) AppleWebKit/537.36 (KHTML  

  \"Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_6) AppleWebKit/537.36 (KHTML  

  \"Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko  

  \"Mozilla/5.0 (Windows NT 6.1; Win64; x64) AppleWebKit/537.36 (KHTML, like  

  \"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 Ge
```

A more secure approach might be a further rotation of proxies between the back and forth sending-receiving process. A proxy server acts as a gateway between the web user and the web server. While the user is exploiting a proxy server, internet traffic flows through the proxy server on its way to the server requested. The request then comes back through that same proxy server and then the proxy server forwards the data received from the website back to the client. The final result will be linear combination of User Agents ID and Proxy server for each sending requests, grating a high security level. Many proxy servers are offered in a paid version, so in this case since security barriers are not that high they will not be implemented. As a further disclaimer many online services are providing free proxies server access, but this comes at a personal security cost due to a couple of reasons: - Free plan Proxies are shared among a number of different clients, so as long as someone has used them in the past for illegal purposes the client is indirectly inheriting their legal infringements. - Very cheap proxies, for sure all of the ones free, have the activity redirected on their servers monitored, profiling in some cases a user privacy violation issue.

2.3.2 Handlers and Trycatches

During scraping many difficulties are met. Some of them might come from website structure issues, so that rooted-tree hierarchies are changed as a consequence of a restructuring. Some others might interest content architecture where data is reallocated to some other places in the webpage, then CSS query are no more able to catch data. Handlers in the form of trycatch error workarounds are explicitly built in this sense. The continuous building and testing of the scraping functioning has required the maintainer to have

a precise and fast debugging experience. The following consideration might give a sense of the time consumed when debugging handlers are not implied: `get.data.catsing()` triggers 34 different scrapping functions that are supposed to point to 34 different data pieces. Within a single function call by default pagination generates 10 pages each of which contains at least 25 different single urls to be scrapped. That leads to a number of 8500 single data information, the probability given 8500 associated to something going missing or misparsed is undoubtedly high. The solution proposed tries to handle fails by implementing as many trycatches as scrapping functions inside the a the single `get.data.catsing()` call. Then inside each single scrapping function are put the handlers. This set up allows to catch (and in some cases prevent) fails starting from the very end of the scraping process. As a consequence of this setting when a scrapping function is not able to gather data an error inside the function is thrown, then the error call is intercepted by the corresponding trycatch, which at the end messages where the error occurs. Below some of the main handlers implied:

- `.get_ua()` verifies that the User Agent in the session 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_~<ef"
  grepl(re, url)
}
```

- `.get_delay()` checks through the robotstxt file if a delay between each request is kindly welcomed. When response is NA delay is not required.

```
.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) ## so that .get_delay results are c
.get_delay(domain = dominio)

## [1] NA
```

2.4 Parallel Computing

Since are opened as many sessions as single links and since for each link are supposed to be called 34 functions run time computation can take a while. Run

time is crucial when dealing with active web pages and time to market in real estate is very important, in here originates the need to have always up-to-date data. Run time optimization involves each level of the scraping process from the “lowest” i.e. inside each single function to the “highest” i.e. the agglomerative function. Inside single scraping functions as a general criteria for loops are avoided due to Rcpp reasons, vectorization is preferred. Within agglomerative function instead the approach was to test two different results. All the following runtime examinations are performed on the `scrape()` functions which is a lightweight version of the final API function. The first attempt was using `furrr` package (Vaughan and Dancho, 2018) which enables mapping through a list with the `purrr`, along with a `future` parallel back end. The approach has shown decent performance, but its run time drastically increases when more requests are sent. This leads to a preventive conclusion about the computational complexity: it has to be at least linear with steep slope. Empirical demonstrations have been made:

Run-Time for First method (furrr multisession)

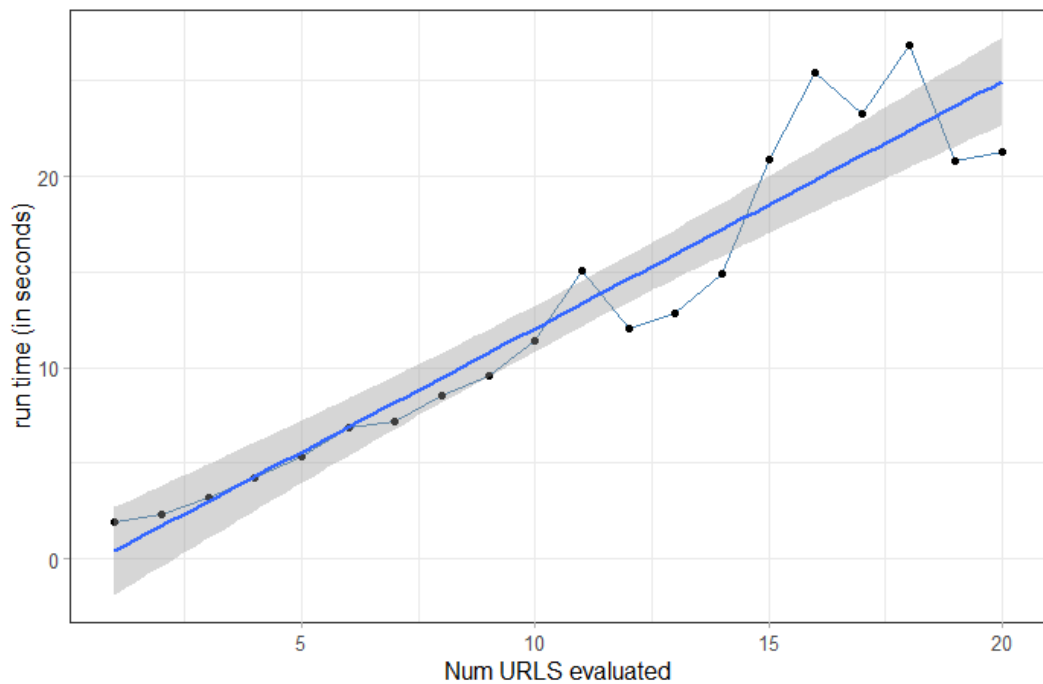


Figure 2.7: computational complexity analysis with Furrr

On the x-axis in the figure the number of urls which are evaluated together, on y axis the run time taken measured in seconds. Iteration after iteration the urls considered are cumulated one at a time. Looking at the blue smoothing curve in between confidence lines the big-O guess might be linear time $\mathcal{O}(n)$, where n are the links considered.

A second attempt tried to explore the **foreach** package (Microsoft and Weston, 2020). This quite recent package enables a new looping construct for executing R code in an iterative way. The core reason for using the **foreach** package is that it supports *parallel execution*, that is, it can execute those repeated operations on multiple processors/cores on the computer, or on multiple nodes of a cluster. The construction follows the r-base looping idea, below steps are summarized:

- start clusters on processors cores
- define the iterator, i.e. “i” equal to the number of elements that are going to be looped
- **.packages**: Inherits the packages that are used in the tasks define below
- **.combine**: Define the combining function that bind results at the end (say `cbind`, `rbind` or `tidyverse::bind_rows`).
- **.errorhandling**: specifies how a task evaluation error should be handle.
- **%dopar%**: the `dopar` keyword suggests `foreach` with parallelization method
- then the function within the elements are iterated
- close clusters

One major concern regards that functions inside the **%dopar%** should be standalone in order to be executed in parallel. For standalone it is meant that everything that is needed to be executed and to output results should be defined inside the **%dopar%**, as it would be opened a new empty environment for each iteration. Moreover as a further consequence packages imported into each clusters, the **.packages** methods takes care of that.

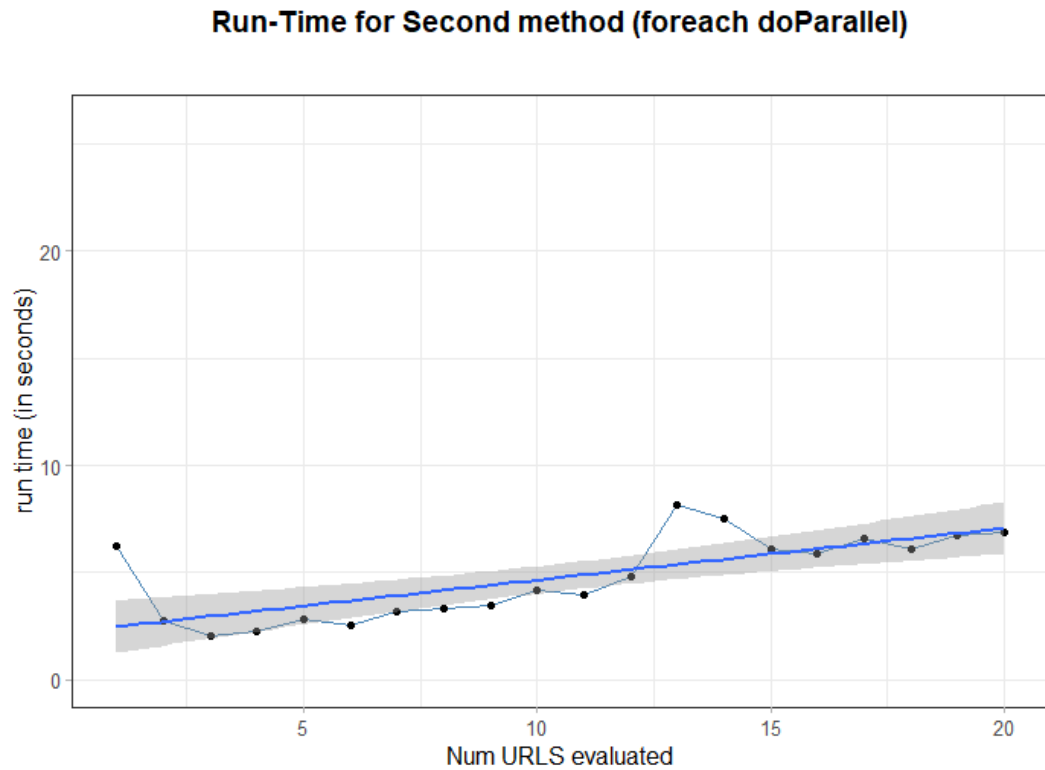


Figure 2.8: computational complexity analysis with Furr

It can be grasped quite easily that the curve now is flattened and a confident guess might be logarithmic time $\mathcal{O}(\log(n))$.

A further performance improvement could be obtained using a new package called `doAzureParallel` which is built on top of the `foreach`. `doAzureParallel` enables different Virtual Machines operating parallel computing throughout Microsoft Azure cloud, but this comes at a substantial monetary cost. This would be a perfect match given that parallel methods seen before accelerates the number of requests sent among different processors or cluster, even though actually the goal is to have something that separates different sessions. Unleashing Virtual Machines allows from one hand to further increase computational capabilities, so the number of potential requests, from the other it can partition requests among different proper machines (a pool of agents for each VM) extending even more the combination of IDs and as a consequence masquerading even better the scraping automation.

2.5 Open Challenges

The main challenge remains unsolved since each single elements have been finely optimized but scraping function and so REST API must be continuously maintained. Indeed What it can not be a-priori optimized are the future changes that involves the website structure. Content architecture as opposite, with some sophistication can take care of finding exact information within the webpage even if the designed is changed. The idea developed in the package Khalil (2018), even though results are not always acceptable, is to crawl the website and to search for targeted keywords. Once keywords are found the scraping algorithm looks for the related information that should be located throughout the html files locally stored. html are known to be very lightweight so computation of this kind is not bothering run time. For this reasons performances with algorithm of this species are very neat but results, as anticipated, are under the expectation. As a solution an accurate text mining approach can be considered as a further enachement. It should be also pointed out that Khalil (2018) is designed to scrape a vast number of websites, as contrary the scraping functions here presented are exclusively designed to be applied on immobiliare.it, even though they can be extended to other related website with no effort. The way the scraping function are designed really facilitates responsive fast debugging but this can not be by any means automatized. The API necessitates frequently to resort to continuous integration (i.e. CI) review to verify the working status. Moreover Error messages can not really be understood sometimes even with handlers, this is due to functions that are called within a parallel beckend that does not allow to print error on console as in this stackoverflow reproducible example¹. So each time an error occurs the “main” functions needs to be taken out of from the parallel back end and separately evaluated. This is time consuming but for the time being no solutions have been found.

¹<https://stackoverflow.com/questions/10903787/how-can-i-print-when-using-dopar>

2.6 Legal Profiles (ancora non validato)

“Data that is online and public is always available for all” is never a good answer to the question “Can I use those web data to my scope”. Immobiliare.it² is not providing any open source data from its own database neither the perception is that it is planning to do so in the future. Immobiliare has not even provided a paid API through which data might be accessed. A careful reading of their terms, reviewed with a intellectual property expert, has been done to get this service running without any legal consequence, as a reference the full policy can be seen in their specialized section³. Nevertheless the golden standard for scraping was respected since the `robot.txt` is neat allowing any actions as demonstrated above. So if it might be the case of misinterpretation of their policy, it will be also the case of lack of communication between servers response and immobiliare.it intent to preserve their own intellectual property. What it was shockingly surprising are the low barriers to obtain information with respect to other counterpart online players. Best practices are applied and delayed requests (even though not asked) have been sent to normalize traffic congestion. But scraping criteria followed are once again fully based on common shared best practises (see section 2.2), and *not* any sort of general agreements between parties. As a result a plausible approach could be applying scraping procedures without any prevention. It would not surely cause any sort of disservice for the website since budget constraints are set low, but in the long run it will cause lagging as soon as budget or subjects will increase. Totally different was the approach proposed by Idealista.com, which is a comparable to immobiliare.it. Idealista does block requests if they are not in compliance with their servers inner rules. User agents in this case must be rotated quite frequently and as soon as a request does not fall within the pool of user agents (i.e. is labeled as web bot) it is immediately blocked and 404 response is sent back. Delay is kindly asked and it must be specified, consequently this slows down scraping function per se.

²<https://www.immobiliare.it/>

³<https://www.immobiliare.it/terms/>

- Idealista content is composed by Javascript so and html parser can not get that.
- Idealista blocks also certain web browser that have a demonstrated “career” in scraping procedures.

All of this leads to accept that entry barriers to scrape are for sure higher than the one faced for Immobiliare. The reticence to share data could be a reflex on how big idealista is; as a matter of fact it has a heavy market presence in some of the Europe real estate country as Spain and France. So the hidden intention was to raise awareness on scraping procedure that in a certain remote way can hurt their business. This has been validated by the fact that prior filtering houses on their website a checkbox has to be signed. The checkbox make the user sign an agreement on their platform according to which data can not be misused and it belongs their intellectual property.

Chapter 3

REST API Infrastructure

In order to provide a fast and easy to use API service to the end user many technologies have been involved. Challenges in scraping as pointed out in section 2.5 are many and still some remains unsolved. Challenges regards not only scraping per se, but also the way the service has to interact with the users. Interactions are many and consequently are problems arised. API service has to be fast otherwise data become obsolete and so happen to the analysis that have relied on those data. Service has to be deployed so that it can be shared over a wider range of clients. Service has to be scalable at need since, due to deployment, when the number of users increases the run time performance should not decrease. Moreover from one hand service has to be continuously integrated and reviewed so that each function can be responsive to immobiliare.it changes. But on the other code behind the service has to be version controlled and freezed, so that when packages are updated service can prevent fails. API has to be also secured granting access only to the ones authorized. In the end Service has to be run at a certain given times and storing data on cloud database, so that it can be tracked back the evolution of the phenomenon under inspection. Open source solutions for each of the requirements stated are available for back-end and front-end integration. Moreover documentations related to technologies served are able to offer flexible solutions to be embedded into the R ecosystem. For all the

requirements the idea is to provide a REST Plumber API with 4 endpoints which calls parallelized scraping functions built in section 2. On top of that a daily Cron Job scheduler, exposing one API endpoint, produces and later stores a .csv file in a NOSQL mongoDB Atlas cloud database. Containerization happens through a Linux OS (Ubuntu distr) Docker container hosted by a AWS EC2 server. API endpoints are secured with https protocols and protected with authentication by nginx reverse proxy. On a second server a Shiny App calls one endpoint with specified parameters which returns daily data from the former infrastructure.

Technologies involved are:

- GitHub version control
- Scheduler cron job, section 3.1
- Docker containers, section 3.2
- Plumber REST API, section 3.3.1
- NGINX reverse proxy, section 3.4
- AWS (Amazon Web Services) EC2 3.5
- MongoDB Atlas
- Shiny, see chapter ??

As a side note each single part of this thesis has been made according to the same API inspiring criteria of reproducibility and self containerization. RMarkdown (Allaire et al., 2020) documents (book's chapters) are compiled and then converted into .html files. Through Bookdown (Xie, 2016) the resulting documents are put together according to general .yaml instruction file and are readable as gitbook. Files are then pushed to a Github repository¹. By a simple trick with GH pages, .html files are displayed into a Github subdomain hosted at link². The resulting deployed gitbook can also produce a .pdf version output through a Xelatex engine. Xelatex compiles .Rmd documents according to a .tex template which formatting rules are contained in a further .yaml

¹<https://github.com/NiccoloSalvini/thesis>

²<https://niccolosalvini.github.io/thesis/>

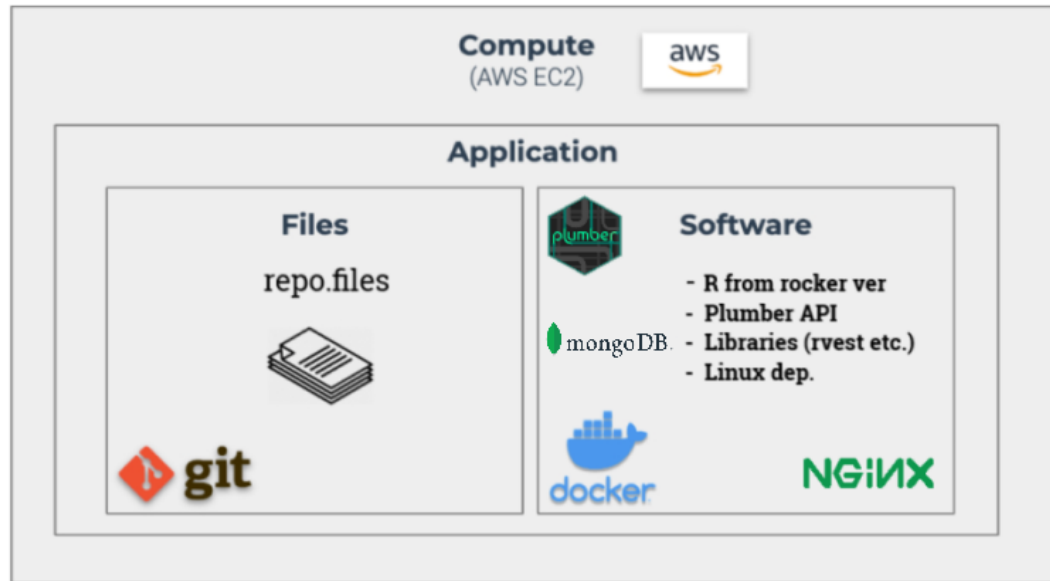


Figure 3.1: complete infrastructure (Matt Dancho source)

file. The pdf version of the thesis can be obtained by clicking the download button, then choosing pdf output version in the upper banner. For further references on the topic Xie (2016)

Some of the main technologies implied will be viewed singularly, nonetheless for brevity reasons some needs to be skipped.

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. In 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 refereces, such as the scheduling algorithm the reader can refer to (Wikiversità, 2020).

3.1.1 Cron Jobs

Cron job is a software utility which acts as a time-based job scheduler in Unix-like OS. Linux users that set up and maintain software environments exploit cron to schedule their day-to-day routines to run periodically at fixed times, dates, or intervals. It typically automates system maintenance but its usage is very flexible to whichever needed. It is lightweight and it is widely used since it is a common option for Linux users. The tasks by cron are driven by a crontab file, which is a configuration file that specifies a set of 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.

Each line of a crontab file represents a job, and has this structure

Each line of a crontab file represents a job. This example runs a shell named scheduler.sh at 23:45 (11:45 PM) every Saturday. .sh commands can update mails and other minor routines.

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

Some rather unusual scheduling definitions for crontabs can be found in this reference (Wikipedia contributors, 2020). Crontab's syntax completion can be made easier through this³ GUI.

³<https://crontab.guru/>

```

# _____ 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 3.2: crontab

The cron job needs to be ran on scraping fuctions at 11:30 PM every single day. The `get_data.R` script first sources an endpoint function, then it applies the function with fixed parameters. Parameters describe the url specification, so that each time the scheduler runs the `get_data.R` collects data from the same source. Day after day .json files are generated and then stored into a NOSQL *mongoDB* database whose credentials are public. Data are collected on a daily basis with the explicit aim to track day-by-day changes both in the new entries an goners in rental market, and to investigate the evolution of price differentials over time. Spatio-Temporal modeling is still quite unexplored, data is saved for future used. Crontab configuration for daily 11:30 PM schedules has this appearance:

```
30 11 * * * /home/oracle/scripts/get_data.R
```

Since now the computational power comes from the machine on which the system is installed. A smarter solution takes care of it by considering run time limits and the substantial inability to share data. To a certain extent what it has been already done since now might fit for personal use: a scheduler can daily execute the scraping scripts and generate a .csv file. Furthermore an application can rely on those data, but evident reasons suggest that it does not suite any need. What it will do the trick would be an open source dedicated software environment or *container* that will contains scraping functions and a scheduler on cloud solving a pair of the problems arisen. This problem can be addressed with a technology that has seen a huge growth in its usage in the

last few years.

3.2 Docker

Docker is a software tool to create and deploy applications using containers. *Docker containers* are a standard unit of software (i.e. software boxes) where everything needed for applications, such as libraries or dependencies can be run reliably and quickly. Furthermore they are also portable, in the sense that they can be taken from one computing environment to the following. Docker containers by default run on kernel Linux OS. Containers can be thought as an abstraction at the app layers that groups code and dependencies together. One major advantage of containers is that multiple containers can run on the same machine with the same OS. Each container can run its own isolated process in the user space, so that each task is complementary to the other. Containers are lightweight and take up less space than Virtual Machines (container images are files which can take up typically tens of MBs in size), can handle more applications and require fewer Virtual Machines and OSs.

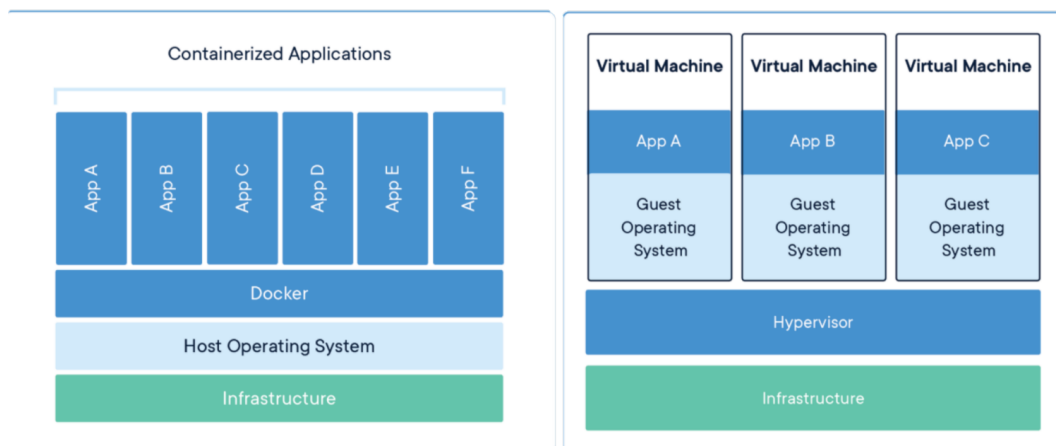


Figure 3.3: docker container vs VM

When containers are built *Docker container Images* are created and can be open sourced through Docker Hub. *Docker Hub* is a web service provided by Docker for searching and sharing container images with other teams or

developers in the community. Docker Hub can connect with GitHub behind authorization entailing an image version control tool. Once the connection is established changes that are pushed with git to the GitHub repository are passed to Docker Hub. The push command automatically triggers the image building. Then docker image can be tagged (salvini/api-immobiliare:latest) so that on one hand it is recognizable and on the other can be reused in the future. Once the building stage is completed the DH repository can be pulled and then run locally on machine or cloud, see section 3.5. Docker building and testing images can be very time consuming. R packages can take a long time to install because code has to be compiled, especially if using R on a Linux server or in a Docker container. Rstudio package manager⁴ includes beta support for pre-compiled R packages that can be installed faster. This dramatically reduces packages time installation (Nolis, 2020). In addition to that an open source project named rocker⁵ has narrowed the path for developers by building custom R docker images for a wide range of usages. What can be read from their own website about the project is: “The rocker project provides a collection of containers suited for different needs. find a base image to extend or images with popular software and optimized libraries pre-installed. Get the latest version or a reproducible fixed environment.”

3.2.1 Why Docker

Indeed⁶, an employment-related search engine, released an article on 2019 displaying changing trends from 2015 to 2019 in Technology Job market. Many changes are relevant in key technologies. Two among the others technologies (i.e. docker and Azure) have experienced a huge growth and both refer to the same demand input: *containers*. 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 reports drive better decisions In

⁴<https://packagemanager.rstudio.com/client/#/>

⁵<https://www.rocker-project.org/images/>

⁶<https://it.indeed.com/>

2020 - Businesses need apps to empower better decision making at all levels

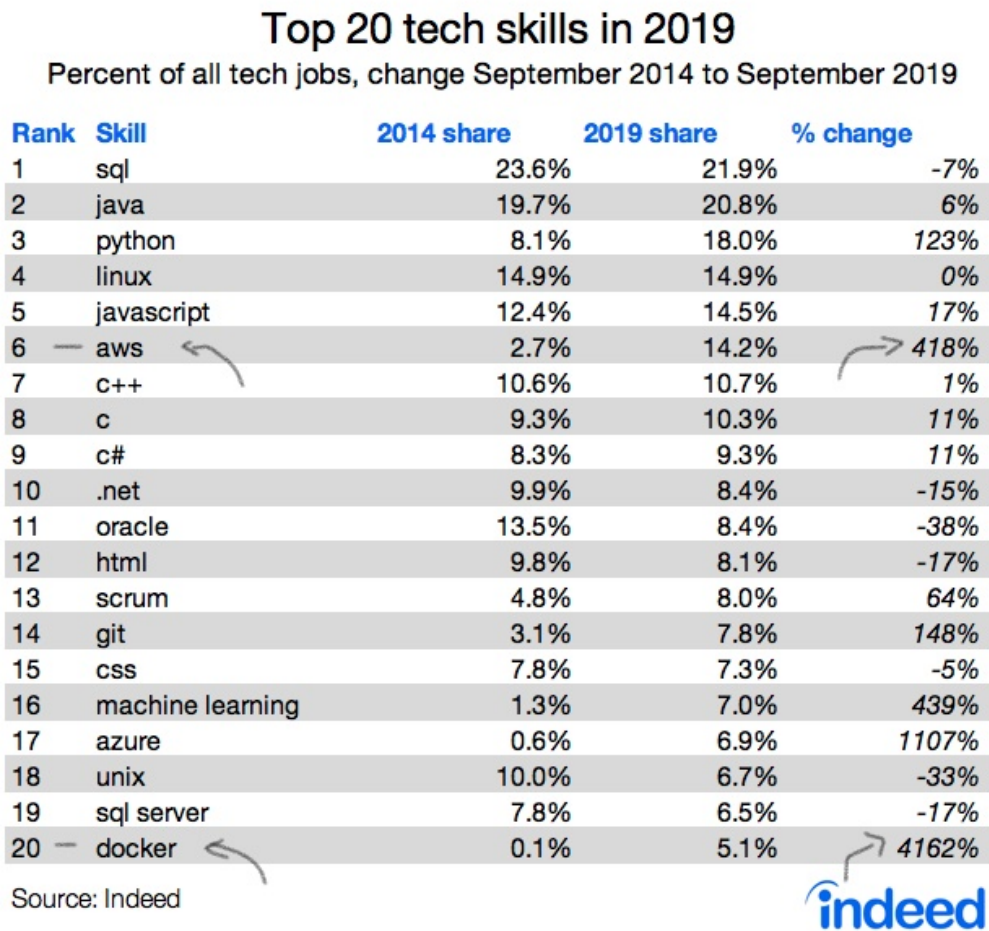


Figure 3.4: docker-stats

For all the things said what docker is bringing to business (Inc., 2020b):

- *Speed 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 transferred 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. In addition due to Docker Hub it is possible to establish a connection between Git and DockerHub. Version
- *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.
- *Fault isolation* : Docker reduces effort and risk of problems with application dependencies. Docker also freezes the environment to the preferred packages version so that it guarantees continuity in deployment and isolate the container from system fails coming from package version updates.

The way to tell docker which system requirements are needed in the newly born software is a *Dockerfile*.

3.2.2 Dockerfile

Docker can build images automatically by reading instructions from a Dockerfile. A Dockerfile is a text 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 sequentially several command-line instructions. For thesis purposes a dockerfile is written with the specific instructions and then the file is pushed to GitHub repository. Once pushed DockerHub automatically parses the repository looking for a plain text file whose name is “Dockerfile”. When It is matched then it triggers the building of the image.

The Dockerfile used to trigger the building of the service docker container has the following set of instructions:

```
FROM rocker/tidyverse:latest

MAINTAINER Niccolo Salvini "niccolo.salvini27@gmail.com"

RUN apt-get update && apt-get install -y \
    libxml2-dev \
    libudunits2-dev

# install R packages
RUN R -e "install.packages(c('magrittr','lubridate', 'plumber', 'rvest', 'stringi', 'jsonlite', 'f

# install 'iterators' dep for DoParallel
RUN R -e "install.packages('https://cran.r-project.org/src/contrib/Archive/iterators/iterators_1.0

# install 'foreach' dep for DoParallel
RUN R -e "install.packages('https://cran.r-project.org/src/contrib/Archive/foreach/foreach_1.4.8.1

# install DoParallel from source since not avail in 4.0.2
RUN R -e "install.packages('https://cran.r-project.org/src/contrib/Archive/doParallel/doParallel_

COPY / /

# expose port
EXPOSE 8000

ENTRYPOINT ["Rscript", "main.R"]
```

Figure 3.5: dockerfile

- `FROM rocker/tidyverse:latest` : the command imports a pre-built image by the rocker team that contains the latest (tag latest) version of base-R along with the tidyverse packages.
- `MAINTAINER Niccolo Salvini "niccolo.salvini27@gmail.com"` : The command tags the maintainer and its e-mail contact information.
- `RUN apt-get update && apt-get install -y \ libxml2-dev \ libudunits2-dev` :The command update and install Linux dependencies needed for running R packages. `rvest` requires `libxml2-dev` and `magrittr` needs `libudunits2-dev`. If they are not installed then associated libraries can not be loaded. Linux dependencies needed have been found by trial and error while building containers. Building logs messages print errors and suggest which dependency is mandatory.

- `RUN R -e "install.packages(c('plumber','tibble','...'),dependencies=TRUE)`
: the command install all the packages required to execute the files (R files) containerized for the scraping. Since all the packages have their direct R dependencies the option `dependencies=TRUE` is needed.
- `RUN R -e "install.packages('https://cran.r-project.org/.../iterators, type='source') RUN R -e "install.packages('https://cran.r-project.org/.../ type='source') RUN R -e "install.packages('https://cran.r-project.org/.../ type='source')` DoParallel was not available in package manager for R version later than 4.0.0. For this reason the choice was to install a previous source version by the online repository, as well as its dependencies.
- `COPY \` The command tells Docker copies all the files in the container.
- `EXPOSE 8000` : the commands 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 (Inc., 2020a))
- `ENTRYPOINT ["Rscript", "main.R"]` : the command tells docker to execute the file `main.R` within the container that triggers the API start. In `main.R` it are pecified both the port and the host where API expects to be exposed (in this case port 8000).

In order to make the system stand-alone and make the service available to a wider range of subjects a choice has to be made. The service has to have both the characteristics to be run on demand and to specify query parameters.

3.3 REST API

API stands for application programming interface and it is a set of definitions and protocols for building and integrating application software. 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 impacting positively on the budget due to resource savings. APIs are thought of as contracts, with documentation that represents an general agreement between parties. There are many types of API that exploit different media and architectures to communicate with apps or services. The specification REST stands for REpresentational State Transfer and is a set of architectural principles. When a request is made through a REST API it transfers a representation of the state to the requester. This representation, is submitted in one out of the many available formats via HTTP: JSON (Javascript Object Notation), HTML, XLT, TXT. JSON is the most popular because it is language agnostic (wha, 2018), as well as more comfortable to be read and parsed. In order for an API to be considered RESTful, it has to conform to these criteria:

(rivedi elenco) - A client-server architecture made up of clients, servers, and resources, with requests managed through HTTP. - Stateless client-server communication, meaning no client information is stored between requests and each request is separate and unconnected. - Cacheable data that streamlines client-server interactions. - A uniform interface between components so that information is transferred in a standard form. This requires that: - resources requested are identifiable and separate from the representations sent to the client. - resources can be manipulated by the client via the representation they receive because the representation contains enough information to do so. - self-descriptive messages returned to the client have enough information to describe how the client should process it. - hypermedia, meaning that after accessing a resource the client should be able to use hyperlinks to find all other currently available actions they can take. - A layered system that organizes each type of server (those responsible for security, load-balancing, etc.) involved the retrieval of requested information into hierarchies, invisible to the client.

REST API accepts http requests as input and elaborates them through end points. An end point identifies the operation through traditional http methods (e.g. /GET /POST) that the API caller wants to perform. Further documentation and differences between HTTP and REST API can be found to this reference⁷.

open REST API examples: - BigQuery API: A data platform for customers to create, manage, share and query data. - YouTube Data API v3: The YouTube Data API v3 is an API that provides access to YouTube data, such as videos, playlists, and channels. - Cloud Natural Language API: Provides natural language understanding technologies, such as sentiment analysis, entity recognition, entity sentiment analysis, and other text annotations, to developers. - Skyscanner Flight Search API: The Skyscanner API lets you search for flights & get flight prices from Skyscanner's database of prices, as well as get live quotes directly from ticketing agencies. - Openweathermap API: current weather data for any location on Earth including over 200,000 cities.

3.3.1 Plumber REST API

Plumber allows the user to create a REST API by adding decoration comments to the existing R code. Decorations are a special type of comments that suggests to Plumber where and when the API specifications parts are. Below a simple example extracted by the documentation:

```
# plumber.R

## Echo back the input
## @param msg The message to echo
## @get /echo
function(msg="") {
```

⁷https://docs.aws.amazon.com/it_it/apigateway/latest/developerguide/http-api-vs-rest.html

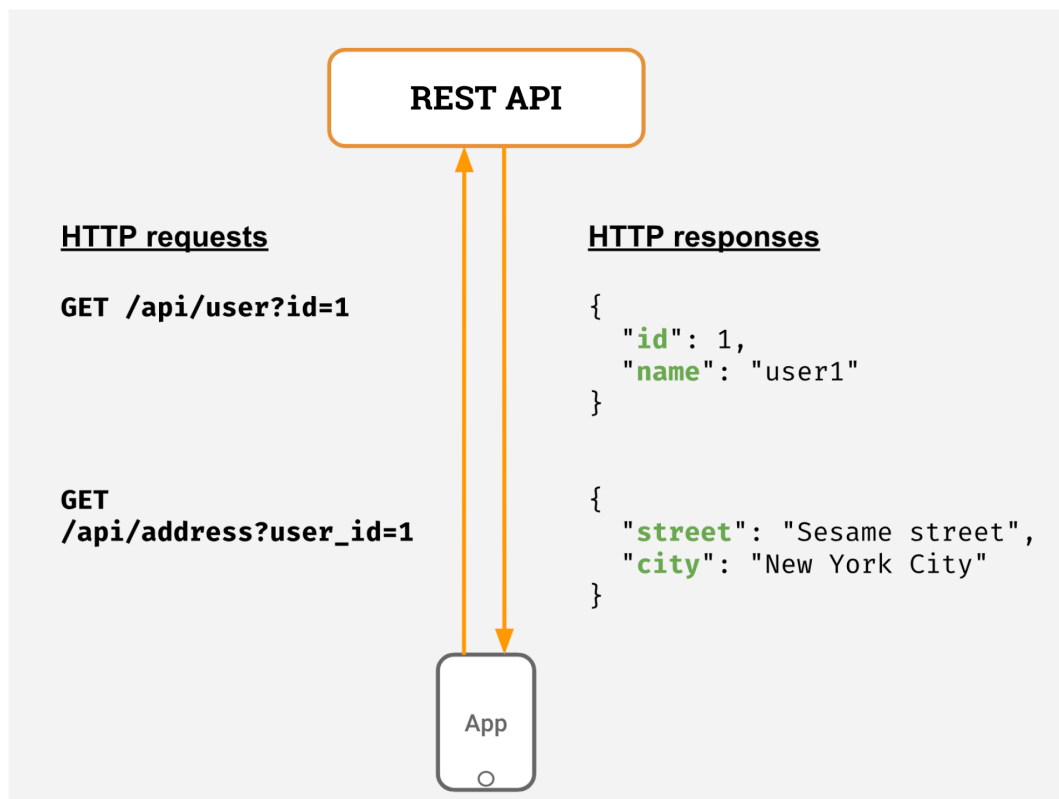


Figure 3.6: API functioning

```

    list(msg = paste0("The message is: '", msg, "'"))
}

## Plot a histogram
## @serializer png
## @get /plot
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)
}

```

three endpoints associated to 2 /GET and 1 /POST requests are made available. Functions are made clear without names so that whenever the endpoint is called functions are directly executed. Decorations are marked as this *##* and they are followed by specific keywords denoted with *@*. - the *@params* keyword refers to parameter that specifies the corpus of the HTTP request, i.e. the inputs with respect to the expected output. If default parameters are inputted then the API response is the elaboration of the functions with default parameters. As opposite endpoint function elaborates the provided parameters and returns a response. - *## @serializer* specifies the extension of the output file when needed. - *## @get* specifies the method of HTTP request sent. - */echo* is the end point name. - *@filter* decorations activates a filter layer which are used to track logs and to parse request before passing the argbody to the end

points.

Many more options are available to customize plumber API but are beyond the scope, a valuable resource for further insights can be found in the dedicated package website (?)

3.3.2 Immobiliare.it REST API

The API service is composed by 4 endpoints */scrape* , */links*, */complete* and */get_data*:

- **/scrape* performs a fast scraping of the website that leverages a shortest path to directly extract 5 covariates from url. url from which data extraction takes place might be composed through parameters. By default the end point scrape data from Milan real estate rental market. Fast scraping is reached thanks to avoiding to access to single links. It is a superficial scraping and does not contain geospatial, however it might fit for regression settings.
- **/links*: extracts the list of single links belonging to each of the page, looking at section 2.1.1 each 25 single links for each sibling. It displays sufficient performances in terms of run time. It is propaedeutic to apply the following endpoint.
- **/complete*: both the function *all.links* and *complete* are sourced. The former with the aim to grab each single links and store it into an object. The latter to actually iterate scraping on each of the links.
- **/get_data*: it triggers the data extraction by sourcing the */complete* endpoint and then storing .json file into the NOSQL mongoDB ATLAS

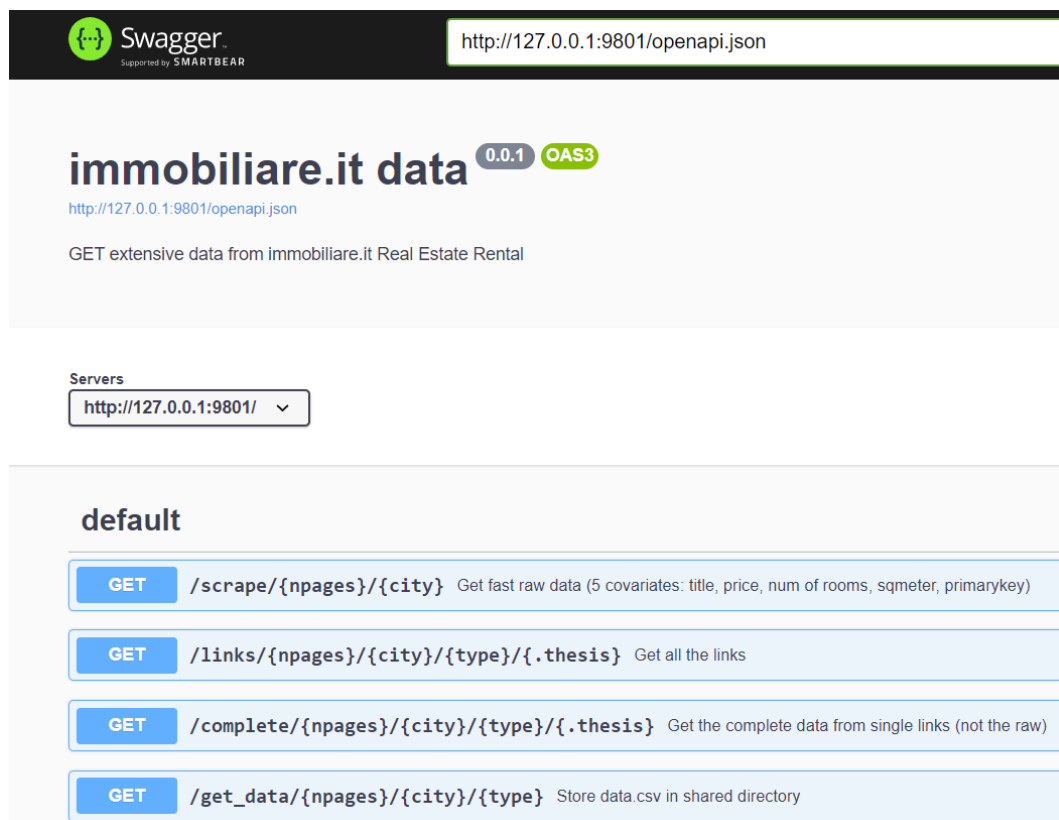


Figure 3.7: swagger

3.3.3 REST API documentation

- Get FAST data, it covers 5 covariates: title, price, num of rooms, sqmeter, primarykey

```
GET */scrape
```

```
@param city [chr string] the city you are interested in (e.g. "roma", "m
```

```
@param npages [positive integer] number of pages to scrape, default = 10
```

```
@param type [chr string] "affitto" = rents, "vendita" = sell
```

```
@param macrozone [chr string] avail: Roma, Firenze, Milano, Torino; e.g.
```

```
content-type: application/json
```

- Get all the links

```
GET */link
```

```
@param city [chr string] the city you are interested to extract data (lo
```

```
@param npages [positive integer] number of pages to scrape default = 10,
```

```
@param type [chr string] "affitto" = rents, "vendita" = sell
```

```
@param .thesis [logical] data used for master thesis
```

```
content-type: application/json
```

- Get the complete set of covariates (52) from each single links, takes a while

```
GET */complete
```

```
@param city [chr string] the city you are interested to extract data (lo
```

```
@param npages [positive integer] number of pages to scrape default = 10,
```

```
@param type [chr string] "affitto" = rents, "vendita" = sell
```

```
@param .thesis [logical] data used for master thesis
```

```
content-type: application/json
```

3.4 NGINX reverse proxy server

For analysis purposes NGINX is open source software for reverse proxying and load balancing. Proxying is typically used to distribute the load among several servers, seamlessly show content from different websites, or pass requests for processing to application servers over protocols other than HTTP. [...]

When NGINX proxies a request, it sends the request to a specified proxied server, fetches the response, and sends it back to the client. It is possible to proxy requests to an HTTP server (another NGINX server or any other server) or a non-HTTP server (which can run an application developed with a specific framework, such as PHP or Python) using a specified protocol. Supported protocols include FastCGI, uwsgi, SCGI, and memcached. [...]

.conf file and installation on Linux server. Security and Authentication.

3.5 AWS EC2 server

Executing REST API on a public server allows to share data with a various number of services thorough multitude of subjects. Since it can not be specified a-priori how many times and users are going to enjoy the service a scalable solutio might fill the needs. Scalable infrastructure through a flexible cloud provider combined with nginx load balancing can offer a stable and reliable infrastructure for a relatively cheap price. AWS offers a wide range of services each of which for a wide range of budgets and integration. Free tier servers can be rent up to a certain amount of storage and computation that nearly 0s the total bill. The cloud provider also has a dedicated webpage to configure the service needed with respect to the usage named amazon cost manager⁸. Amazon Elastic Compute Cloud (EC2) is a web service that contributes to a secure, flexible computing capacity in the AWS cloud. EC2 allows to rent as many virtual servers as needed with customized capacity, security and storage.

⁸<https://aws.amazon.com/en/aws-cost-management/>

[few words still on EC2]

3.5.1 Launch an EC2 instance

The preliminary step is to pick up an AMI (Amazon Machine Image). AWS AMI are already-set-up machines with stadardized specification designed to speed up the process of choosing the a customed machine. Since the project is planned to be nearly 0-cost a “Free Tier Eligible” server is chosen. By checking the Free Tier box all the available free tiers are displayed. The machine selected has this specification: t2.micro with 1 CPU and 1GB RAM and runs on a Ubuntu distribution OS. First set up settings needs to be left as-is, net-working and VPC can always be updated when needed. In the “add storage” step 30 GB storage are selected, moreover 30 represent the upper limit since the server can be considered free tier. Tags windows are beyond the scope. Secondly configuration needs to account security and a new entry below SSH connection (port 22) has to be set in. New security configuration has to have TCP specification and should be associated to port 8000. Port 8000, as in dockerfile section 3.2.2, has been exposed and needs to be linked to the security port opened.

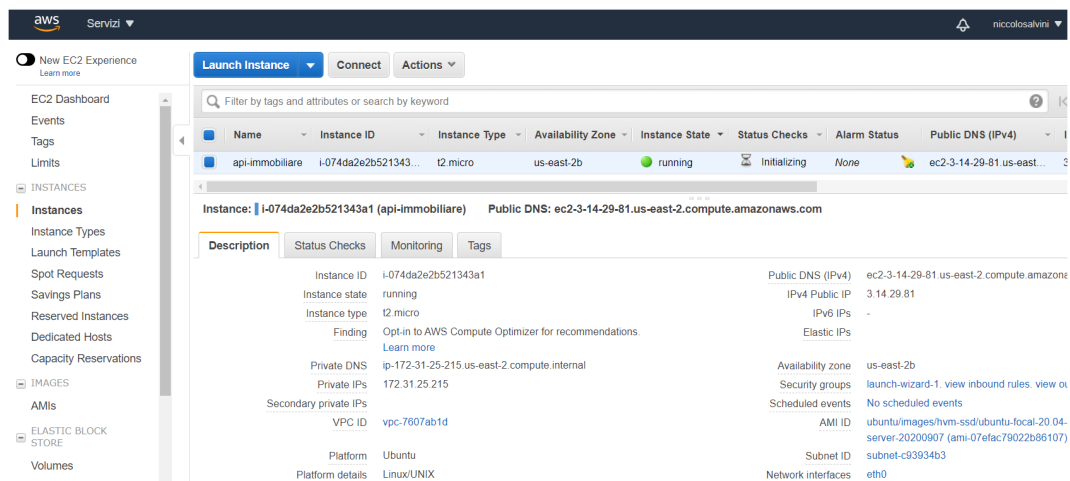


Figure 3.8: aws_dashboard

At this point instance is prepared to run and in a few minutes is deployed. Key

pairs, if never done before, are generated and .pem file is saved and securely stored. Key pairs are mandatory to access to the Ubuntu server via SSH. SSH connection in Windows OS can be handled with PuTTY⁹, which is a SSH and telnet client designed for Windows. At first PuTTYgen has to convert the key pair .pem file into a .ppk extension (otherwise Putty can not read it). Once converted .ppk is sourced in the authorization panel. If everything works and authentication is verified then the Ubuntu server CLI appears and an interaction with the server is made available.

3.6 Further Integrations

Pins is an r packages this link¹⁰

⁹<https://www.putty.org/>

¹⁰https://rstudio.com/resources/rstudioconf-2020/deploying-end-to-end-data-science-with-shiny-plumber-and-pins/?mkt_tok=eyJpIjoiTmprNU1USXhPVEprWXpNMSIsInQiOiJtTUhKVzlvSjVIV2hKc0NRNVU1NTRQYSsrRGd5MWMy

Chapter 4

INLA computation

INLA (Rue et al., 2009) stands for Integrated Nested Laplace approximation and constitutes a computational alternative to traditional MCMC methods. INLA does approximate bayesian inference on special type of models called LGM (Latent Gaussian Models) due to the fact that they are *computationally* convenient. The benefits are many, some among the other are:

- Low computational costs, even for large models.
- It provides high accuracy.
- Can define very complex models within that framework.
- Most important statistical models are LGM.
- Very good support for spatial models.
- Implementation of spatio-temporal model enabled.

INLA uses a combination of analytics approximations and numerical integration to obtain an approximated posterior distribution of the parameters in a shorter time period. The chronologic steps in the methodology presentation follows the course sailed by Moraga (2019) blended with the author choice to skip details. As a matter of fact the aim of the chapter is to provide a comprehensive intuition oriented to the immediate application of the methodology, without stepping too long on mathematical details. By the way details e.g model assessment and control options are handled under the hood by the

package and can be tuned within the main function, most of them are covered by Gómez Rubio (2020). Notation is imported from Marta Blangiardo (2015) and Gómez Rubio (2020), and quite differ from the one presented in the original paper by Rue, Chopin and Martino (2009). As further notation remarks: bold symbols are considered as vectors, so each time they occur they have to be considered like the *ensemble* of their values. In addition $\tilde{\pi}$ in section 4.2 are the Laplace approximation of the underlying integrals. Moreover the inner functioning of Laplace approximation and its special usage within the INLA setting is far from the scope, but an easy shortcut oriented to INLA is in Marta Blangiardo (2015).

INLA can fit only Latent Gaussian type of models and the following work tries to encapsulate its properties. Then afterwards a problem can be reshaped into the LGM framework with the explicit purpose to explore its benefits. When models are reduced to LGMs then joint posterior distribution can be rewritten and then approximated with INLA. A hierarchical bayesian structure on the model will help to integrate many parameter and hyperparameter levels and simplify interpretation. Generic Information on the project and the R-INLA package are contained in the introduction to last section 4. In the end a brief application on a toy spatial dataset is proposed with the aim to fasten the familiarity with the methodology and to come to grip with INLA results.

4.1 Latent Gaussian Models LGM

Given some observations $y_{i...n}$ in order to define a Latent Gaussain Model within the bayesian framework it is convenient to specify at first an *exponential family* (Gaussian, Poisson, Exponential...) distribution function characterized by some parameters ϕ_i (usually expressed by the mean $E(y_i)$) and some other hyper-parameters $\psi_k, \forall k \in 1 \dots K$. The parameter ϕ_i can be defined as an additive *latent linear predictor* η_i , as pointed out by Krainski and Rubio ((2019)) through a link function $g(\cdot)$, i.e. $g(\phi_i) = \eta_i$. A comprehensive

expression of the linear predictor takes into account all the possible effects on covariates

$$\eta_i = \beta_0 + \sum_{m=1}^M \beta_m x_{mi} + \sum_{l=1}^L f_l(z_{li})$$

where β_0 is the intercept, $\beta = \{\beta_1, \dots, \beta_M\}$ are the coefficient that quantifies the linear effects on covariates $x = (x_1, \dots, x_M)$ and $f_l(\cdot), \forall l \in 1 \dots L$ are a set of random effects defined in terms of a z set of covariates $z = (z_1, \dots, z_L)$ (e.g. `rw`, `ar1`). As a consequence of the last assumption the class of LGM can receive a wide range of models e.g. GLM, GAM, GLMM, linear models and spatio-temporal models. This constitutes one of the main advantages of INLA, which can fit many different models, starting from simpler and ending with more complex. Contributors recently are extending the methodology to many areas as well as models moreover they are trying to incorporate INLA with non gaussian latent models as Rubio (2020) pointed out. All the latent components can be conveniently grouped into a variable denoted with θ such that: $\theta = \{\beta_0, \beta, f\}$ and the same can be done for hyper parameters $\psi = \{\psi_1, \dots, \psi_K\}$. Then the probability distribution conditioned to parameters and hyper parameters is then:

$$y_i \mid \theta, \psi \sim \pi(y_i \mid \theta, \psi)$$

Since data (y_1, \dots, y_n) is drawn by the same distribution family but it is conditioned to parameters which are conditional independent (i.e. $\pi(\theta_i, \theta_j \mid \theta_{-i,j}) = \pi(\theta_i \mid \theta_{-i,j}) \pi(\theta_j \mid \theta_{-i,j})$) (Rue and Held, 2005) then the joint distribution is given by the product of all the independent parameters i.e. the likelihood. Moreover the Product operator index i ranges from 1 to n , i.e. $\mathbf{I} = \{1 \dots n\}$. When an observation is missing so the corresponding $i \notin \mathbf{I}$ INLA automatically will not include it in the model avoiding errors (2020). As a consequence the likelihood expression is:

$$\pi(y \mid \theta, \psi) = \prod_{i \in \mathbb{I}} \pi(y_i \mid \theta_i, \psi) \quad (4.1)$$

Each data point is connected to one combination θ_i out of all the possible linear combinations of elements in θ *latent field*. The latent aspect of the field regards the undergoing existence of many parameter combination alternatives. Furthermore hyper parameters are by definition independent, in other words ψ will be the product of many univariate priors (Gómez Rubio, 2020). A Multivariate Normal distribution is imposed on the latent field θ such that it is centered in 0 with precision matrix $Q(\psi)$ (the inverse of the covariance matrix $Q^{-1}(\psi)$) depending only on ψ hyper parameter vector i.e., $\theta \sim \text{Normal}(\mathbf{0}, Q^{-1}(\psi))$. As a notation remark some authors choose to keep the covariance matrix expression as Q and its inverse precision matrix as Q^{-1} . This is strongly not encouraged for two reasons: the first is that the default hyperparameter option in INLA R package uses the precision matrix, the second it over complicates notation when writing down conditional expectation as Rue pointed out *miss lit*. However notation for covariance function introduced in chapter 5.2.1 i.e. Matérn has to be expressed through covariance matrix, this passage will be cleared out in the dedicated section so that confusion is avoided. The exponential family density function is then expressed through:

$$\pi(\theta \mid \psi) = (2\pi)^{-n/2} |Q(\psi)|^{1/2} \exp\left(-\frac{1}{2} \theta Q(\psi) \theta\right) \quad (4.2)$$

The conditional independence assumption on the latent field θ leads $Q(\psi)$ to be a sparse precision matrix since for a general pair of combinations θ_i and θ_j the resulting element in the precision matrix is 0 i.e. $\theta_i \perp \theta_j \mid \theta_{-i,j} \iff Q_{ij}(\psi) = 0$ (2015). A probability distribution function with those characteristics is said *Gaussian Markov random field* (**GMRF**). GMRF as a matter of fact are Gaussian variables with Markov properties which are encoded in the precision matrix Q (Rue et al., 2009). (puoi dire di più) From here it comes the source of run time computation saving, inherited using GMRF for inference. As a

consequence of GMRF representation of the latent field, matrices are sparse so numerical methods can be exploited (Marta Blangiardo, 2015). *Moreover when Gaussian Process (see chapter 5.1), which are used to integrate spatial components, are represented as GMRF through SPDE (Stochastic Partial Differential Equations) approach, then INLA can be used as a computing choice. This last assumption will be framed in chapter 4 and verified in chapter 6.* Once priors distributions are specified for ψ then the joint posterior distribution for θ and ψ is

$$\pi(\theta, \psi | y) \propto \underbrace{\pi(\psi)}_{\text{prior}} \times \underbrace{\pi(\theta | \psi)}_{\text{GMRF}} \times \underbrace{\prod_{i=1}^n \pi(y_i | \theta_i, \psi)}_{\text{likelihood}}$$

Last expression is said a Latent Gaussian Models, **LGM**, if the whole set of assumptions imposed since now are met. Therefore all models that can be reduced to a LGM representation are able to host INLA methodology. Then plugging in the *likelihood* (4.1) and *GMRF* (4.2) expression the posterior distribution can be rewritten as

$$\begin{aligned} \pi(\theta, \psi | y) &\propto \pi(\psi) \times \pi(\theta | \psi) \times \pi(y | \theta, \psi) \\ &\propto \pi(\psi) \times \pi(\theta | \psi) \times \prod_{i=1}^n \pi(y_i | \theta_i, \psi) \\ &\propto \pi(\psi) \times |Q(\psi)|^{1/2} \exp\left(-\frac{1}{2}\theta'Q(\psi)\theta\right) \times \prod_i^n \exp(\log(\pi(y_i | \theta_i, \psi))) \end{aligned}$$

And by joining exponents by their multiplicative property it is obtained

$$\pi(\theta, \psi | y) \propto \pi(\psi) \times |Q(\psi)|^{1/2} \exp\left(-\frac{1}{2}\theta'Q(\psi)\theta + \sum^n \log(\pi(y_i | \theta_i, \psi))\right) \quad (4.3)$$

4.2 Approximation in INLA setting

INLA is not going to try to estimate the whole posterior distribution from expression (4.3). Instead it will try to estimate the posterior marginal distribution effects for each θ_i combination in the latent parameter θ , given the hyper parameter priors specification ψ_k . Proper estimation methods however are beyond the scope of the analysis, further excellent references are suggested in their respective part by Rubio (2020) in section 2.2.2 and Blangiardo & Cameletti (2015) in section 4.7.2. The marginal posterior distribution function for each latent parameter element θ_i is

$$\pi(\theta_i | y) = \int \pi(\theta, \psi | \mathbf{y}) \pi(\psi | \mathbf{y}) d\psi \quad (4.4)$$

The posterior marginal integral for each hyper parameter ψ_k , $k = 1, \dots, K$ is

$$\pi(\psi_k | y) = \int \pi(\psi | y) d\psi_{-k}$$

where the notation ψ_{-k} is a vector of hyper parameters ψ without considering k th element ψ_k .

The goal is to have approximated solution for latent parameter posterior distributions. To this purpose A *hierarchical procedure* is now imposed since the “lower” hyper parameter integral, whose approximation for the moment does not exist, is nested inside the “upper” parameter integral that takes hyper param as integrand. Hierarchical structures are welcomed very warmly since they are convenient later in order to fit a hierarchical bayesian model approached in the next chapter ??). While many approximation strategies are provided and many others are emerging for both the hyper param and for the latent field, the common ground remains to unnest the structure in two steps such that:

- step 1: compute the Laplace approximation $\tilde{\pi}(\psi | y)$ for each hyper

parameters marginal: $\tilde{\pi}(\psi_k | y)$

- step 2: compute Laplace approximation $\tilde{\pi}(\theta_i | \psi, y)$ marginals for the parameters given the hyper parameter approximation in step 1: $\tilde{\pi}(\theta_i | y) \approx$

$$\int \tilde{\pi}(\theta_i | \psi, y) \underbrace{\tilde{\pi}(\psi | y)}_{\text{Estim. in step 1}} d\psi$$

Then plugging approximation in the integral observed in (4.4) it is obtained:

$$\tilde{\pi}(\theta_i | y) \approx \int \tilde{\pi}(\theta_i | \psi, y) \tilde{\pi}(\psi | y) d\psi$$

In the end INLA by its default approximation strategy through *simplified Laplace approximation* uses the following numerical approximation to compute marginals:

$$\tilde{\pi}(\theta_i | y) \approx \sum_j \tilde{\pi}(\theta_i | \psi^{(j)}, y) \tilde{\pi}(\psi^{(j)} | y) \Delta_j$$

where $\{\psi^{(j)}\}$ are a set of values of the hyper param ψ grid used for numerical integration, each of which associated to a specific weight Δ_j . The more the weight Δ_j is heavy the more the integration point is relevant. Details on how INLA finds those points is beyond the scope, but the strategy and grids seraches are offered in the appendix follwing both Rubio and Blangiardo.

4.2.1 further approximations (prolly do not note include)

INLA focus on this specific integration points by setting up a regular grid about the posterior mode of ψ with CCD (central composite design) centered in the mode (Gómez Rubio, 2020).

The approximation $\tilde{\pi}(\theta_i | y)$ can take different forms and be computed in different ways. Rue et al. (2009) also discuss how this approximation should be in order to reduce the numerical error (Krainski, 2019).

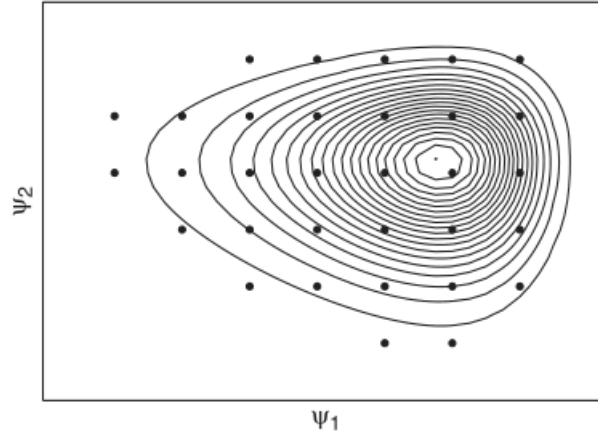


Figure 4.1: CCD to spdetoy dataset, source Marta Blangiardo (2015)

Following Gómez Rubio (2020), approximations of the joint posterior for the hyper parameter $\tilde{\pi}(\psi_k | y)$ is used to compute the marginals for the latent effects and hyper parameters in this way:

$$\tilde{\pi}(\psi | \mathbf{y}) \propto \frac{\pi(\theta, \psi, y)}{\tilde{\pi}_G(\theta | \psi, y)} \Big|_{\theta=\theta^*(\psi)}$$

In the previous equation $\tilde{\pi}_G(\theta | \psi, y)$ is a gaussian approximation to the full condition of the latent effect $\theta^*(\psi)$ is the mode for a given value of the hyper param vector ψ

At this point there exists three types of approximations for $\pi(\theta | \psi, y)$

- first with a gaussian approximation, estimating mean $\mu_i(\psi)$ and variance $\sigma_i^2(\psi)$.
- second using the *Laplace Approximation*.
- third using *simplified Laplace Approximation*

(rivedere meglio)

4.3 R-INLA package in a bayesian hierarchical regression perspective

4.3.1 Overview

INLA computations and methodology is developed by the R-INLA project whose package is available on their website¹. Download is not on CRAN (the Comprehensive R Archive Network) so a special source repo link, which is maintained by authors and collaborators, has to be optioned. The website offers also a forum where a daily discussion group is opened and an active community is keen to answer. Moreover It also contains a number of reference books, among which some of them are fully open sourced as gitbook. Furthermore as Havaard Rue has pointed out in a web-lecture on the topic, the project is gaining importance due to its new applications and recent use cases, but by no means it is replacing the older MCMC methods, rather INLA can integrate pre existing procedures. The core function of the package is `inla()` and it works as many other regression functions like `glm()`, `lm()` or `gam()`. Inla function takes as arguments the formula (where are response and linear predictor), the data (expects a data.frame obj) on which estimation is desired together with the distribution of the data. Many other methods inside the function can be added through lists, such as `control.family` and `control.fixed` which let the analyst specifying priors distribution both for θ parameters, ψ hyper parameters and the outcome precision τ , default values are non-informative. `control.fixed` as said regulates prior specification through a plain list when there only a single fixed effect, instead it does it with nested lists when fixed effects are greater than 2, a guided example might better display the behaviour: `control.fixed = list(mean = list(a = 1, b = 2, default = 0))` In the chunk above it is assigned prior mean equal to 1 for fixed effect “a” and equal 2 for “b”; the rest of the prior means are set equal to 0. Inla objects are `inla.dataframe` summary-type lists containing

¹<http://www.r-inla.org>

the results from model fitting. Results contained in the object are specified in the table below, even though some of them requires special method: (se riesco più elegante in kable) Following Krainski & Rubio (2019) observations $y(s_1), \dots, y(s_n)$ are taken from a toy generated dataset and a hierarchical linear regression is fitted.

Function	Description
<code>summary.fixed</code>	Summary of fixed effects.
<code>marginals.fixed</code>	List of marginals of fixed effects.
<code>summary.random</code>	Summary of random effects.
<code>marginals.random</code>	List of marginals of random effects.
<code>summary.hyperpar</code>	Summary of hyperparameters.
<code>marginals.hyperpar</code>	List of marginals of the hyperparameters.
<code>mlik</code>	Marginal log-likelihood.
<code>summary.linear.predictor</code>	Summary of linear predictors.
<code>marginals.linear.predictor</code>	List of marginals of linear predictors.
<code>summary.fitted.values</code>	Summary of fitted values.
<code>marginals.fitted.values</code>	List of marginals of fitted values.

Figure 4.2: summary table list object, source: Krainski (2019)

4.3.2 Linear Predictor

SPDEtoy dataset, that has a spatial component, is generated from a y_i Gaussian variable; its moments are μ_i and precision τ .

The formula that describe the linear predictor has to be called directly inside the `inla()` function or it can be stored in the environment into a variable. The mean moment in the gaussian distribution μ_i is expressed as the *linear predictor* η_i (i.e. $E(y_i | \beta_0, \dots, \beta_M, x_{i1}, \dots, x_{iM}) = \eta_i$). The function that maps the linear predictor into the parameter space is identity as in section ?? i.e. $\eta_i = \beta_0 + \sum_{m=1}^M \beta_m x_{mi} + \sum_{l=1}^L f_l(z_{li})$. After including s_1 and s_2 spatial covariates the linear predictor takes the following form: $\beta_0 + \beta_1 s_{1i} + \beta_2 s_{2i}$, where once again β_0 is the fixed effect i.e. intercept and the β_j are the linear effect on covariates. INLA allows also to include non-linear effects with the

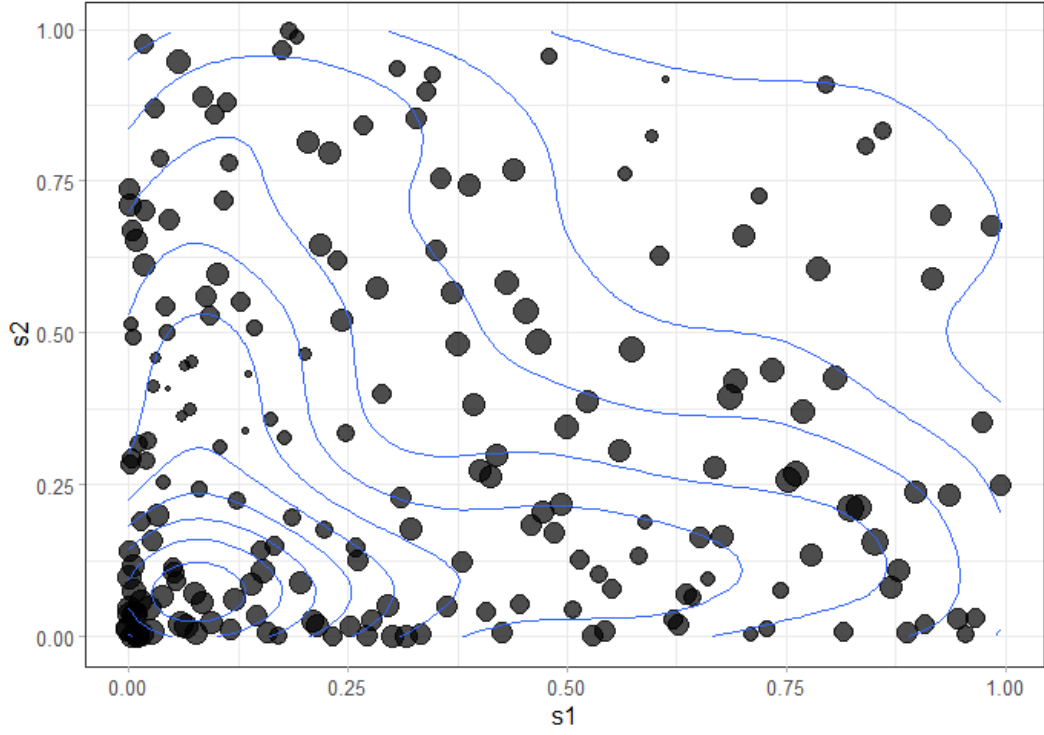


Figure 4.3: SPDEtoy plot, author's source

$\mathbf{f}()$ method inside the formula. \mathbf{f} are fundamental since they are used to incorporate the spatial component in the model through the Matérn covariance function, this will be shown in section (boh). Once the formula is decided then priors has to be picked up; for the intercept a customary choice is uniform. Prior for Gaussian latent parameters are vague and they have 0 mean and 0.001 precision, then the prior for τ is a Gamma with parameters 1 and 0.00005. Prior initial choice can be later adapted.

The summary of the model parameters is:

$$y_i \sim N(\mu_i, \tau^{-1}), i = 1, \dots, 200$$

$$\mu_i = \beta_0 + \beta_1 s_{1i} + \beta_2 s_{2i}$$

$$\beta_0 \sim \text{Uniform}$$

$$\beta_j \sim N(0, 0.001^{-1}), j = 1, 2$$

$$\tau \sim Ga(1, 0.00005)$$

```
data("SPDEtoy")
formula = y ~ s1 + s2
m0 = inla(formula, data = SPDEtoy)
```

```
\begin{table}
```

```
\caption{(\#tab:table_INLA)prova}
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode	
(Intercept)	10.1321487	0.2422118	9.6561033	10.1321422	10.6077866	10.1321497	7
s1	0.7624296	0.4293757	-0.0814701	0.7624179	1.6056053	0.7624315	7
s2	-1.5836768	0.4293757	-2.4275704	-1.5836906	-0.7404955	-1.5836811	7

```
\end{table}
```

The output offers among the others a summary of the posterior marginal values for intercept, coefficient and covariates, as well as precision. Below the plots for the parameters and hyperparameters. From the summary it can be seen that the mean for s2 is negative, so the more the value of the y-coordinates increases the more the output decreases, that is truer looking at the SPDEtoy cotour plot. Plots can be generated by calling the `plot` function on the `inla` object, however the one generated below are `ggplot2` outputs coming from the `$marginals.fixed` list object.

```
\begin{table}
```

```
\caption{(\#tab:higer_posterior_density_interval)higer posterior density interval}
```

	low	high
level:0.9	-2.291268	-0.879445

```
\end{table}
```

R-Inla also has r-base fashion function to compute statistics on marginal posterior distributions for the density, distribution as well as the quantile

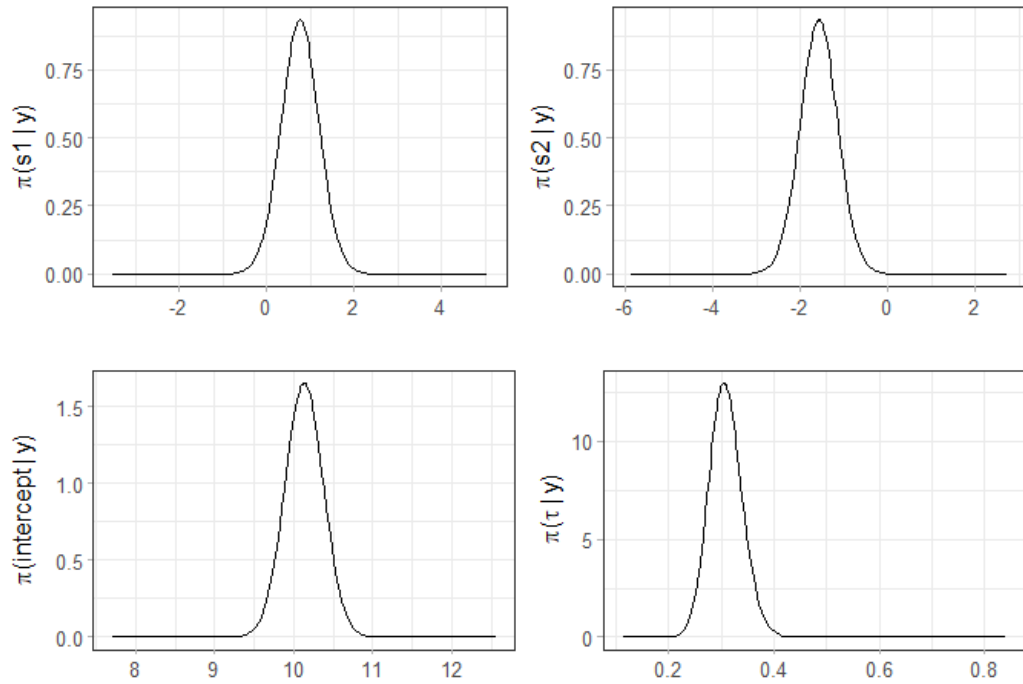


Figure 4.4: linear predictor marginals, author's creation

function respectively `inla.dmarginal`, `inla.pmarginal` and `inla.qmarginal`. One major option which is conveniently packed into a dedicated function computes the higher posterior density credibility interval `inla.hpdmarginal` for a given covariate's coefficient, such that

$$\int_{q_1}^{q_2} \tilde{\pi}(\beta_2 | y) d\beta_2 = 0.90 \text{ with } .1 \text{ Confidence Level.}$$

Recall that the interpretation is different from the frequentist: in Bayesian statistics β_j comes from probability distribution, while frequentists considers β_j as fixed unknown quantity whose estimator (random variable conditioned to data) is used to infer the value (2015).

Chapter 5

Point Referenced Data Modeling

Geostatistical data are a collection of samples of geo type data indexed by coordinates (e.g. latlong, eastings and northings) that originate from a spatially continuous phenomenon (Moraga, 2019). Data as such can monitor a vast range of phenomena, as an example disease cancer detection (Bell et al., 2006) at several sites, COVID19 spread in China (Li et al., 2020), PM pollution concentration in a North-Italian region Piemonte (Cameletti et al., 2012). Moreover house prices variation, as observed in Gómez Rubio (2020), where selling prices smoothly vary between closer neighborhoods. All the Examples taken before might document a spatial nature of data according to which closer observations can display similar values, this phenomenon is named spatial autocorrelation. Spatial autocorrelation conceptually originates from geographer Waldo Tobler whose famous quote, known as first law of geography, inspires geostatisticians:

“Everything is related to everything else, but near things are more related than distant things”

— Waldo R. Tobler

Spatial models are explicitly designed to take into account this behavior and

can separate spatial patterns from simply random spatial variance. Spatial data can be partitioned into three spatial data type whose modeling tools are specific with respect to their category.

- Areal Data
- **Point Referenced Data**
- Point Pattern Data

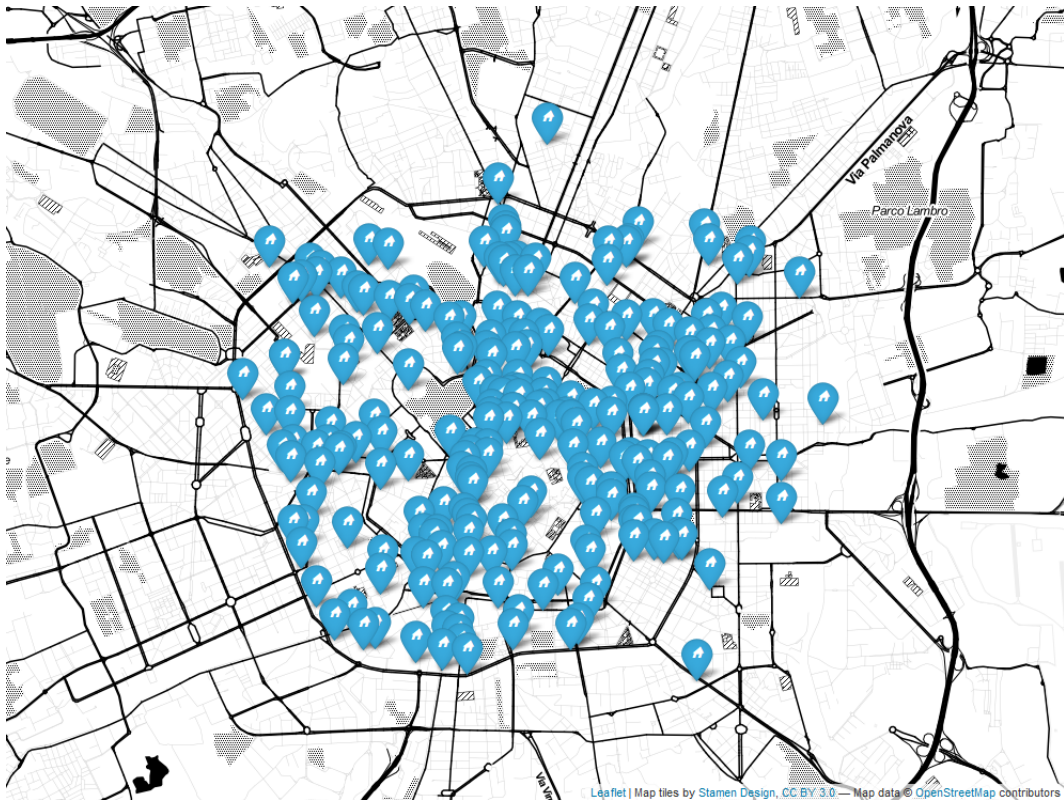


Figure 5.1: point referenced data example, Milan Rental Real Estate, Author's Source

REST API seen in chapter 3 extracts point referenced data, so modeling methodologies described in this analysis will exclusively take into account point referenced oriented techniques. In order to extend the notion from discrete measurements (i.e. point referenced) to a continuous spatial surface a stochastic process, namely Gaussian Process, has to be introduced and constrained according to convenient properties. GP are then evaluated with a specific covariance function, i.e. Matérn. The reason why Matérn is selected

as candidate for covariance function will be much more clear in the next chapter 6. Hedonic Price Models are at first introduced and then a brief literature review is offered. Hedonic Prices brings to this work the theoretical basis but they do not suggest estimation methods, which are essentially the major issue in geostatistics. For this reason Hedonic Models are exploited into a spatial bayesian regression framework with the aim to apply INLA (seen in chapter 4) methodology. At first standard Bayesian regression is presented as introduction, then the spatial component in the form of a GP is added to the model. Many parameters are considered so far, as a consequence a hierarchy structure is imposed. To this extent an entire section is dedicated to hierarchy which simplifies model building and methodology understanding as well as allowing to bring in many different parameters that come from different levels through the exchangeability property. As a matter of fact parameters originate from the Gaussian latent field, but also from Matérn covariance function tuning hyper parameters. Then INLA is applied and a GMRF representation of GP is... Spatial kriging is essential to predict the process at new locations so that the spatial surface can be plotted and analyzed. In the end models have to be checked and verified with resampling schemes which are once again specific to the data type and the scope of the analysis.

(forse mettere alla fine come further developments) As a side note Spatial data can also be measured according to a further dimension which is the Time. Latest literature suggests that spatio temporal models are the most accurate, as a consequence it might be interesting to research time correlation between subsequent spatial data time points, a valuable reference is offered in Paci et al. (2017). This will not take an enormous effort due to the fact that on a daily basis REST API generates data which are stored as .json file on a DB. Future research on this data might consider the idea to include the time component in the model.

5.1 Gaussian Process (GP)

For simplicity let's consider y point of interest observations

$y(s_1), y(s_2), \dots, y(s_n)$ from a random spatial process Y , such that:

$Y(s_1), Y(s_2), \dots, Y(s_n)$ observed at location s_1, \dots, s_n . In the context of geostatistical data each observation has to be considered as a partial realization of an unobserved random spatial process. $\{Y(s) : s \in D \subset \mathbb{R}^2\}$, where surface D is a subset of r -dimensional Euclidean space \mathbb{R}^r . Moreover

When $r = 1$ it is the most simple stochastic process widely explored in literature i.e. time series process. However geostatistical data always have $r = 2$ (i.e. lat and long, eastings and northings) or eventually $r = 3$, when elevation data is available. The stochastic process Y is observed in a fixed set of “monitoring stations” and inference can be done regarding moments of the realized process. This information are essential to build a spatially continuous surface over the y -studied variable in order to predict the phenomenon at locations not yet observed.

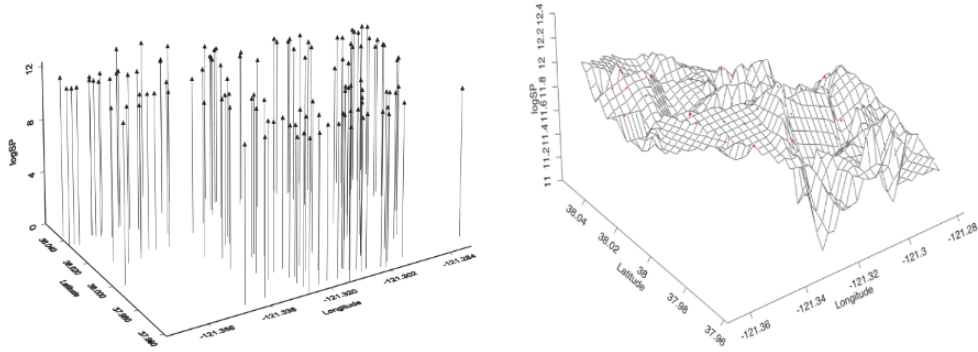


Figure 5.2: 3D scatterplot and surface, Stockton data.

Definition 5.1 (GP definition). A collection of n random variables, such as $Y(s_1), Y(s_2), \dots, Y(s_n)$ that are *valid* spatial processes are said to be a **GP**

if for any set of spatial index n and for each set of corresponding locations $\{y(s_1), \dots, y(s_n)\}$ follows a multivariate *Gaussian* distribution with mean $\mu = \{\mu(s_1), \dots, \mu(s_n)\}$ and covariance matrix $\mathbf{Q}_{i,j}^{-1}, \forall i \neq j$

Even though sometimes it is more convenient to express the covariance matrix as its inverse i.e. precision matrix $Q_{i,j}$ (Marta Blangiardo, 2015). The covariance matrix relates each observation to each of the others through a covariance function defined as $\mathcal{C}(\cdot)$.

GP in the spatial context must check two important properties in order to exploit INLA, even though both of these assumptions can be relaxed:

- **Stationary.**
- **Isotropy.**

Stationarity in a stochastic process can be *strong*, *weak* or *intrinsic*. The strong property forces the distribution of the process $\{y(s_1), \dots, y(s_n)\}$ for any given spatial index n and its correspondent location sets $s_{1,\dots,n}$ to be the same as the one in $\{y(s_1 + h), \dots, y(s_n + h)\}$, where h is a number

belonging to R^2 . On the other hand the weak property ensures that if the

GP mean moment is constant over the study domain $\mu(\mathbf{s}) \equiv \mu$ (e.g. $E[Y(s)] = \mu, \forall s \in D$) then the covariance functions does depend only on the distance (euclidean $\|s_i - s_j\|$ distance) between each couple points.

Weak stationarity consequences are the most interesting: It does not matter whether observations are placed either in a specific region, nor the direction towards they are oriented, the covariance functions $\mathcal{C}(h)$ can summarize the

process through the separation vector \mathbf{h} i.e. $\mathcal{C}(\mathbf{s}, \mathbf{s} + \mathbf{h}) = \mathcal{C}(\mathbf{h}), \forall \mathbf{h} \in \mathbb{R}^r$ (Banerjee et al., 2014). In other words weak stationarity in GP implies being

invariant under *translation* (2019). The relationship between strong and weak is not bijective since being strong implies also being weak, but the

opposite is not always true for non-Gaussian process. Furthermore through the intrinsic stationary property it is meant that $E[Y(\mathbf{s} + \mathbf{h}) - Y(\mathbf{s})] = 0$, the

second moment of the latter expression can be written as

$E[Y(\mathbf{s} + \mathbf{h}) - Y(\mathbf{s})]^2$ leading to $\text{Var}(Y(\mathbf{s} + \mathbf{h}) - Y(\mathbf{s}))$. Last expression is called *variogram* and can be expressed with $2\gamma(\mathbf{h})$, even though its half, i.e. $\gamma(\mathbf{h})$, is more interpretable, namely *semivariogram* (Cressie, 2015).

Semivariograms are characterized by mainly 3 tuning parameters:

- *range* σ^2 : At some offset distance, the variogram values will stop changing and reach a sort of “plateau”. The distance at which the effect occurs is called the range $\frac{\Delta\gamma(\mathbf{h})}{h} \approx 0$.
- *sill* τ^2 : The “plateau” value at which the variogram stops changing $\frac{\Delta\gamma(\mathbf{h})}{h} = 0$.
- *nugget* $\tau^2 + \sigma^2$: The discontinuity at the origin. Although this theoretically should be zero, sampling error and short scale variability can cause it to be non-zero $\gamma(0)$.

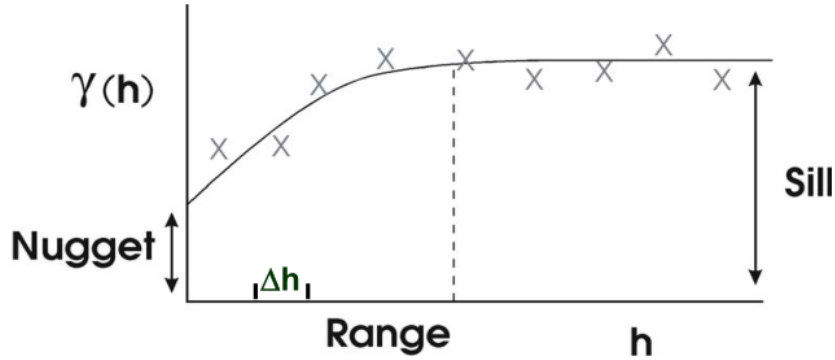


Figure 5.3: variogram example

presi i dati con le relative distanze euclidee a coppie di punti si binnano le distanze grazie ad un offset ottenendo i valori per il semivariogram. ottenuti i valori si fitta il semivariogram a quei valori, un modo è la likelihood. A questo punto si calcolano le tre grandezze nugget sill e range per poi poter far uscire le funzioni di covarianza.

The process is said to be **Isotropic** if the covariance function depends only on the between-points distance $\|\mathbf{h}\|$ so it is invariant under *rotation* (2019). A

further way of seeing the property is that Isotropy implies concentric decaying contours that resemble the vanishing of spatial dependence, and so covariance values too. then if the last assumption does not hold and direction towards point are distant from each other matters within the spatial domain D , then is said to be **Anisotropic**. Formalizing the results:

$$\mathcal{C}(\mathbf{h}) = \mathcal{C}(\|\mathbf{h}\|)$$

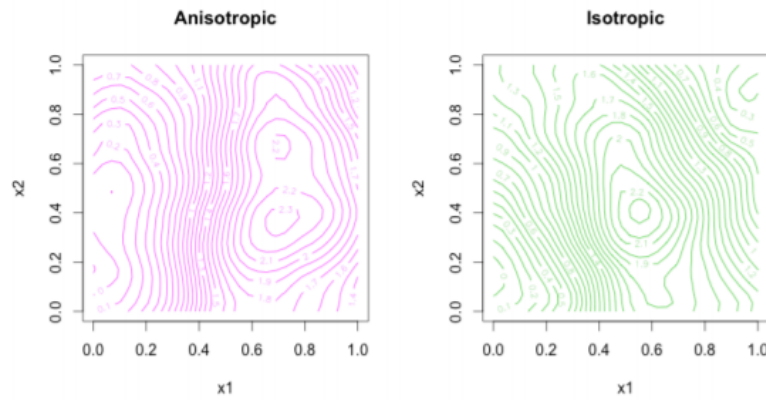


Figure 5.4: isotropy VS anisotropy, source Blanchet-Scalliet et al. (2019)

5.2 Spatial Covariance Function

The covariance function $\mathcal{C}(\cdot)$ ensures that all the values that are close together in input space will produce output values that are close together. $\mathcal{C}(\cdot)$ needs to inherit the *validity* characteristics from the random spatial process, furthermore it has to be *positive definite*. In addition covariance function must share characteristic properties of functions, such as:

(cerca di capire queste...)

- Multiply valid covariance functions (summing independent random variables)
- Mixing covariance functions (mixing distributions)
- Convolving covariance functions, this will be very important ...

Covariance functions under stationary and isotropic GPs displays two important properties: they are constant in mean within D i.e. $\mathcal{C}(\mathbf{s}, \mathbf{s} + \mathbf{h}) = \mathcal{C}(\mathbf{h})$, $\forall \mathbf{h} \in \mathbb{R}^r$ and they depends on distance vector \mathbf{h} , not direction i.e. $\mathcal{C}(\mathbf{h}) = \mathcal{C}(\|\mathbf{h}\|)$ There are many covariance functions and ways to relate distant points on a spatial domain D . Typically the choice of the Covariance can depend either on data or the scope of the analysis. Covariance functions are wrapped into special hyper parameters which are mainly three:

1. *Range*: At some offset distance, the variogram values will stop changing and reach a “plateau”. The distance at which this occurs is called the range.
2. *Sill*: The “plateau” value at which the variogram stops changing.
3. *Nugget*: The discontinuity at the origin. Although this theoretically should be zero, sampling error and short scale variability can cause it to be non-zero

(espressione della covariance function insieme a alle σ^2 come:

$C(\mathbf{s} + \mathbf{h}, \mathbf{s} \mid \theta) = \sigma^2 \mathbf{R}(\|\mathbf{h}\|; \phi)$ spiega anche queste due sotto

$$\mathbf{w} = (w(\mathbf{s}_1), \dots, w(\mathbf{s}_n))' \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{R}(\phi)) \text{ where } \mathbf{R}(\phi)_{ij} = \rho(\|\mathbf{s}_i - \mathbf{s}_j\|; \phi)$$

$$\Sigma_\theta = \sigma^2 \mathbf{R}(\phi) + \tau^2 I_n$$

A summary of the most used covariance functions are presented below.

$$\begin{aligned}
\text{Exponential} \quad \mathcal{C}(\mathbf{h}) &= \begin{cases} \tau^2 + \sigma^2 & \text{if } h = 0 \\ \sigma^2 \exp(-\phi h) & \text{if } h > 0 \end{cases} \\
\text{Gaussian} \quad \mathcal{C}(\mathbf{h}) &= \begin{cases} \tau^2 + \sigma^2 & \text{if } h = 0 \\ \sigma^2 \exp(-\phi^2 h^2) & \text{if } h > 0 \end{cases} \\
\text{Matérn} \quad \mathcal{C}(\mathbf{h}) &= \begin{cases} \tau^2 + \sigma^2 & \text{if } h = 0 \\ \frac{\sigma^2}{2^{\nu-1}\Gamma(\nu)} (\phi h)^\nu K_\nu(\phi h) & \text{if } h > 0 \end{cases}
\end{aligned}$$

5.2.1 Matérn Covariance Function

Matérn is special since when it is used together with a stationary and isotropic GP, the SPDE approach can provide a GMRF representation of the same process, chapter 6 discloses this fundamental property. Matérn can also be accounted as the most used in geostatistics (Krainski et al., 2018) and (Gómez Rubio, 2020) and is tuned mainly by two parameters, a scaling one $\kappa > 0$, usually set equal to the range by the relation $\sigma^2 = \frac{\sqrt{8\lambda}}{\kappa}$ and a smoothing one $\nu > 0$. A *stationary* and *isotropic* Matérn covariance function has this form:

$$\mathcal{C}(\mathbf{h}) = \begin{cases} \tau^2 + \sigma^2 & \text{if } h = 0 \\ \frac{\sigma^2}{2^{\nu-1}\Gamma(\nu)} (\phi t)^\nu K_\nu(\phi t) & \text{if } h > 0 \end{cases}$$

$\Gamma(\nu)$ is a Gamma function depending on ν values, $K_\nu(\cdot)$ is a modified Bessel function of second kind. The smoothness parameter ν in the figure below takes 4 different values showing the potentiality of Matérn to relates distances to covariance values. When $\nu = 1$... When $\nu = 1/2$ it becomes the exponential covariance function, When $\nu = 3/2$ it uncovers a convenient closed form, when $\nu \approx \infty$, in this case for representation purposes $\nu = 80$ it becomes Gaussian covariance function.

ancora di più su matern, forse di più in spde

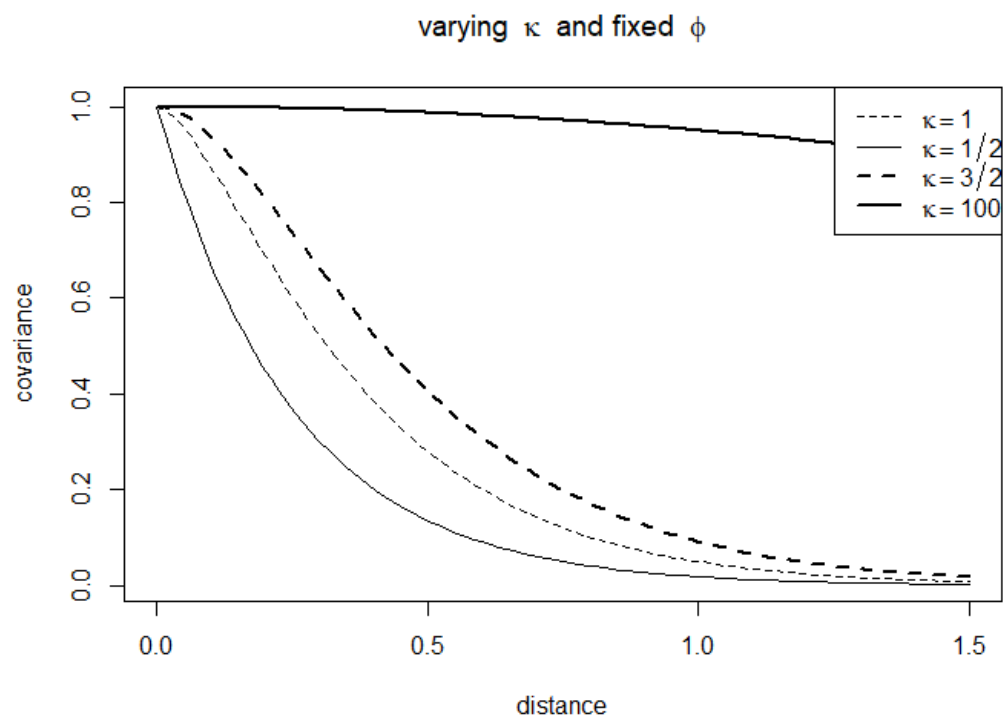


Figure 5.5: matern correlation function for 4 diff values of ν with ϕ fixed, author's source

5.3 Hedonic models Literature Review and Spatial Hedonic Price Models

The theoretical foundation of the Hedonic Price Models (from now on HPM) resides in the consumer utility theory of Lancaster (1966) together with Rosen (1974) market equilibrium. According to Lancaster the utility of a commodity does not exist by itself, instead it exists as the sum of the utilities associated to its separable characteristics. Integrating Lancaster, Rosen introduces HPM and suggests that each separate commodity characteristics are priced by the markets on the basis of supply and demand equilibrium. Applying HPM to Real Estate in a market context, from the buy side house prices (but also rents) are set as the unit cost of each household attributes, conversely from the selling side the expenditures associated to build of each them. Formalizing the results, Hedonic Price P in Real Estate is expressed as a general f functional form that takes as input the house characteristics vector \mathbf{C} .

$$P = f(c_1, c_2, c_3, \dots, c_n)$$

Vector \mathbf{C} since now might contain a unidentified and presumably vast number of ungrouped characteristics. In this setting Malpezzi (2008) tried to organize house features by decomposing \mathbf{C} into mutually exclusive and exhaustive subgroups. An overview of the vector components involved is given by Ling and Ling (2019) according to which P represents the house price, S is the structural characteristics of the house, N represents the neighborhood characteristics, L signifies the locational characteristics, C describes the contract conditions and T is time. β is the vector of the parameters to be estimated. Then

$$P = f(S, N, L, C, T, \beta)$$

Historically a first attempt to include spatial effect in urban economic literature is provided by *Alonso (1964) miss ref.* Its contribution was to raise voice on house prices (also rent) mainly depending on land price and a number of purely spatial covariates like CBD, the distance from City Business District. Other covariates were transport cost per kilometer and community income, even though they were defined also as spatial parameters through distances. The model proposed by Alonso is called monocentric since the centroid from which distances are calculated is only one. Moreover a first touch to spatial data thory was done since the CBD was defined as areal unit with well-defined boundaries of regular or irregular shape. However applications of the model were not convincing since empirical studies offered a different picture. Results instead displayed a Poly-centric areal structure (universities and Malls) which might be better explaining prices. The model also assumed that covariates like CBD are only informative within city center boundaries and then displayed no significance out of the core of the city.

Poly-centric theory was also more coherent with the architectural and socio-economical evolution of cities during that times, therefore mono centric theory was then criticized and abandoned. Critics regarded also neighborhood quality measure and boundary problems *Dubin (1987) miss ref.* Dubin for these reasons developed a model including areal effects in the error term since handling these covariates was posing several hard challenges. Areal data choice for Dubin was forced since he was interested in land values, geostatics interest was not a focus also due to the difficulties in gathering accurate data. Coming to recent literature a change in focus has been made by switching from theory based model to estimation methods. As a consequence to the change in focus Ling and Ling (2019) said that practitioners should spend more time in variable selection and model

specification with respect to their specific need. As Ling has observed the emerging trends are in the field of semi-parametric and non-parametric methods (2019). Historically semi-parametric regression considers models indexed by spatial coordinates *Pace RK (1995)*. At the same time *Kammann and Wand (2003)* gave birth to geoaddivitive models where the spatial component is added as a covariate. [...]

A further aspect of the problem is posed by scholars that do not consider rents to be representative for the actual value of real estate. Nevertheless in empirical analysis rent value are considered a proxy for real estate pricing (Herath and Maier, 2011). A further argument to endorse this hypothesis is brought by Manganelli et al. (2013) considering housing a commodity, then the selling or the rental should be considered interchangeable economic actions with respect to same inner need to be satisfied. This is also truer to the thesis' extent since Manganelli, Morano, and Tajani have centered their analysis exactly on italian real estate data. Moreover Capozza and Seguin (1996) discussed on how much rent-price ratio predicts future changes both in rents and prices. Among all the other discussions raised they brought the decomposition of rent-price ratio into two parts: the predictable part and the unexplained residuals part. The predictable part was discovered to be negatively correlated with price changes, in other words cities in which prices are relatively high with respect to rents are associated with higher capital gains that might justify that misalignment. This is also true for the opposite, that is cities in which prices are lower with respect to the rents, and this effect can not be associated to any local condition, realize lower capital gains. A further argument is offered by Clark (Clark, 1995) which went after the Capozza and Seguin work. Rent-price ratio is negatively correlated with following future changes in rents. In other words prices are still higher when areas in which they are observed documents an increase in rent prices. All the literature review above is oriented to a long-run alignment of price and rent.

5.4 Point Referenced Regression for univariate spatial data

Since in HPM the relationships between the characteristics of the house, i.e. vector \mathbf{C} and the price P is not in any case fixed by econometric literature it is possible to assume any f functional form. The open possibility to apply a wide range of relationship between covariates fit in the INLA setting, since Latent Gaussian Models are prepared to accept a any linear and non linear f functions 4.1 through the `f()` method. Hedonic price models are, as a consequence, a subset of models that can be fitted into LGM and therefore by INLA method.

Moreover what the vast majority of econometric literature (*Greene, 2018*) suggest to apply a is log-linear / square root model. This is due to the fact that log transformation / square root smooths the skewness of prices normalizing the curve, leading to more accurate estimates. Having an exponential family generating process lowers even further computational cost for reasons linked to the $\tilde{\pi}(\psi)$ hyper param INLA approximation (Marta Blangiardo, 2015). Notation is taken from the previous chapter 4, for brevity purposes β \mathbf{X} and y indicates vectors incorporating all their respective realizations and the s spatial component is left out in favor of the observation pedix i .

The simplest log-linear bayesian regression model assumes linear relationship between predictors and a Normal data generating process: (log has been taken out for simplicity, bu it will be then considered in the regression setting) (valuta l'idea che per interpretabilità di modellarla come Gamma exponential family anzichè tenerla normale)

$$\log(y_i) \sim \text{Normal}(\mu_i, \sigma^2)$$

$$y_i = \mu_i + \varepsilon_i$$

then by the following relationship

$E(y_i | \beta_0, \dots, \beta_M, x_{i1}, \dots, x_{iM}) = \beta_0 + \sum_{m=1}^M \beta_m x_{im}$ it is possible to specify a more general linear predictor (seen also in chapter 4) through an identity link

function i.e. $\eta_i = g(\mu_i) = \mu_i$ obtaining:

$$\eta_i = \beta_0 + \sum_{m=1}^M \beta_m x_{mi} + \sum_{l=1}^L f_l(z_{li})$$

Where, once again, the mean structure linearly depends on some \mathbf{X} covariates, β coefficients, $f_l(\cdot), \forall l \in 1 \dots L$ are a set of random effects defined in terms of a z set of covariates $z = (z_1, \dots, z_L)$ (e.g. rw, ar1) and ε_i white noise error. Priors have to be specified and a non informativeness for $\tau^2 = 1/\sigma^2$ and β is chosen, such that $\pi(\tau^2) \propto 1$ and $\pi(\beta) \propto 1$. As a consequence the conditional posterior for the parameters of interest β is:

$$\beta | \sigma^2, y, X \sim \text{MVNormal} \left((X'X)^{-1} X'y, \sigma^2 (X'X)^{-1} \right)$$

where the mean structure corresponds to the OLS estimator: $(X'X)^{-1} X'y$ for β and then to obtain the marginal posterior for β it is needed to integrate with respect to σ^2 .

In order to engage the spatial coordinate components into the regression setting w_i has to be added to the equation. w_i is set as a stationary and isotropic GP with mean 0 and variance as covariance function expressed as Matérn. Recall that GP The new regression setting integrates the *spatial error* part in the name of w_i and a *non-spatial error* part ε_i distributed

normally with mean 0 and variance τ^2 , i.e. $N(0, \tau^2)$, which offers its contribution error to the nuggets via the covariance function. Consequently there is one more parameter to estimate. It is worth mentioning that the distribution of w_i at a finite number of points is considered a realization of a multivariate Gaussian distribution. In this case, the likelihood estimation is possible and it is the multivariate Gaussian distribution with covariance Σ .

$$\log(y_i) = \beta_0 + (\mathbf{X})'\beta + w_i + \varepsilon_i$$

The covariance of the marginal distribution of y_i at a finite number of locations is $\Sigma_y = \Sigma + \tau^2 \mathbf{I}$, where \mathbf{I} denotes the indicator function (i.e., $\mathbf{I}(i = i') = 1$ if $i = i'$, and 0 otherwise). This is a short extension of the basic GF model, and gives one additional parameter to estimate

5.5 Hierarchical Bayesian models

Spatial Models are characterized by many parameters which in turn are tuned by other hyper-parameters. Traditionally Bayesian hierarchical models are not widely adopted since they have high computational burdens, indeed they can handle very complex interactions via random components, especially when dealing with spatio temporal data Ling and Ling (2019). Blangiardo e Cameletti (2015) tried to approach the problem from a different angle offering an intuitive solution on how hierarchy relates different levels parameters. This is done by reversing the problem and starting from data back to parameters, instead the other way round. So taking a few steps back the problem can be reformulated by starting from grouping observation into categories and then trying to impose a hierarchical structure on data based on the categories. As a result observations might fall into different categories, underlining their natural characteristics, such as: some of them might belong to category *levels* like males or females, married or not-married. Moreover

diving into the specific problem house prices can be faceted by which floor they belong or whether they are assigned to different energy classes and many others more. As an example Blangiardo and Cameletti example consider grouping data according to just a single 9 *levels* category. Data for the reasons stated before can be organized such that each single observation (squares in figure below) belongs to its respective mutually exclusive and collectively exhaustive category (circles in figure).

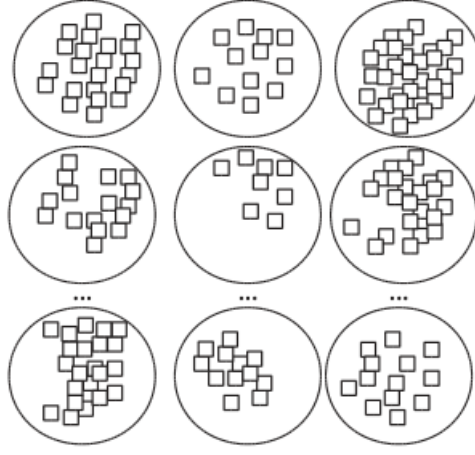
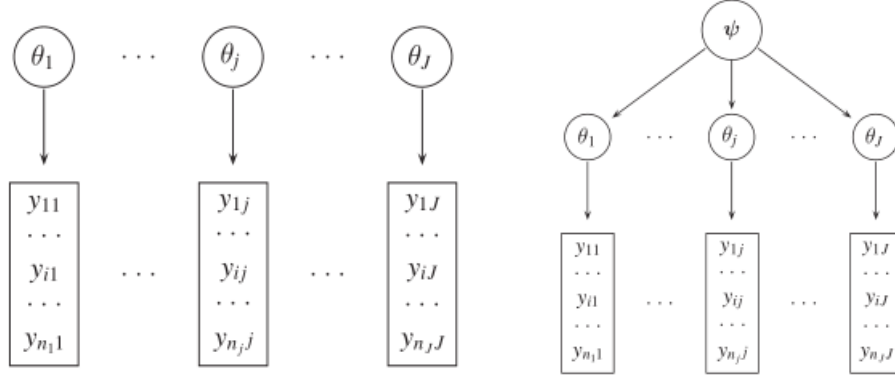


Figure 5.6: 9 levels cat vs observaitions, source Marta Blangiardo (2015)

Furthermore data can be partitioned into two meta-categories, *first level* and *second level*, highlighting the parameter and hyper parameter chain roles.

First level are identified by sampling observations which are drawn by the same probability distribution (squares) . *Second level* (circles) are categories and might be associated to a set of parameters $\theta = \{\theta_1, \dots, \theta_J\}$. Since the structure is hierarchical, a DAG (Directed Acyclical Graph) (2015) representation might sort out ideas. If categories are represented by different θ_j nodes and edges (arrows in the figure) are the logical belonging condition to the category then a single parameter θ model has the right figure form:



To fully take into account the hierarchical structure of the data the model should also consider further levels. Since $\{\theta_1, \dots, \theta_J\}$ are assumed to come from the same distribution $\pi(\theta_j)$, then they are also assumed to be sharing information (Marta Blangiardo, 2015), (left figure). When a further parameter $\psi = \{\psi_1, \dots, \psi_K\}$ is introduced, for which a prior distribution is specified, then the conditional distribution of θ given ψ is:

$$\pi(\theta_1, \dots, \theta_J | \psi) = \int \prod_{j=1}^J \pi(\theta_j | \psi) \pi(\psi) d\psi$$

This is possible thanks to the conditional independence property already encountered in chapter 4, which means that each single θ is conditional independent given ψ . This structure can be extended to allow more than two levels of hierarchy since the marginal prior distributions of θ can be decomposed into the product of their conditional priors distributions given some hyper parameter ψ as well as their prior distribution $\pi(\psi)$.

$$\pi(\theta) = \int \pi(\theta | \psi_1) \pi(\psi_1 | \psi_2) \dots \pi(\psi_{L-1} | \psi_L) \pi(\psi_L) d\psi_1 \dots d\psi_L$$

ψ_l identifies the hyper parameter for the l_{th} level of hierarchy. Each further parameter level ψ is conditioned to its previous in hierarchy level $l - 1$ so that the parameter hierarchy chain is respected and all the linear combinations of

parameters are carefully evaluated. The *Exchangeability* property enables to have higher H nested DAG (i.e. add further L levels) and to extend the dimensions in which the problem is evaluated, considering also time together with space. From a theoretical point of view there are no constraints to how many L levels can be included in the model, but as a drawback the more the model is nested the more it suffers in terms of interpretability and computational power. Empirical studies have suggest that three levels are the desired amount since they offer a good bias vs variance trade-off.

5.6 INLA model through spatial hierarchical regression

INLA model seen in section 4.1 can be rearranged according to the hierarchical structure considering also the regression settings for point referenced data stated in the previous section 5.4.

As an initial step it is required to specify a probability distribution for $y = (y(s_1), \dots, y(s_n)) = (y_1, \dots, y_n)$, this is a mandatory step looking at the 4.3.2 methods needed to compute the `inla()` function. A Normal distribution for simplicity is chosen.

As *first level* is picked up an **exponential family** sampling distribution (i.e. Normally distributed, Gamma one other choice), which is *exchangeable* with respect to the $\theta = \{\beta_0, \beta, f\}$ *latent field* and hyper parameters ψ_1 , which includes also the ones coming from the latent Matérn GP process w_i . The Spatial Guassian Process is centered in 0 and with Matérn covariance function as τ^2 . w_i addresses the spatial autocorrelation between observation through a Matérn covariance function $\mathcal{C}(\cdot|\psi_1)$ which in turn is tuned by hyper param included in ψ_1 . Moreover the w_i surface has to be passed in the formula method definition 4.3.2 via the `f()` function, so that INLA takes into consideration the spatial component.

$$y \mid \theta, \psi_1 \sim N(\beta_0 + (\mathbf{X}_i)' \beta + w_i, \tau^2 I_n) = \prod_{i=1}^n N(y_i \mid \theta_i, \psi_1)$$

Then at the *second level* the latent field θ is characterized by a Normal distribution given the remaining hyper parameters ψ_2 , recall the covariance matrix $Q^{-1}(\psi_2)$, depending on ψ_2 hyperparameters, is handled now by a Matérn covariace function depeding on its hyperparamter. This is done in order to map the GP spatial surface into a GMRF by SPDE solutions.

$$\theta \mid \psi_2 \sim N(0, \mathcal{C}(\cdot, \cdot \mid \psi_2))$$

In the end a *third level* hyper parameters $\psi = \{\psi_1, \psi_2\}$ having some specified prior distribution i.e. $\psi \sim \pi(\psi)$,

5.7 Spatial Kriging

In Geostatistics the main interest resides in the spatial prediction of the spatial latent field pr the response variable at location not yet observed. Assumed the model in the previous section, suppose that y^* is not a observed occurrence of the response variable at location s_0 (not in the data) of the GP w_i spatial surface estimated through observed refereced points in y . As a consequence of exchangeability (first step previous section 5.6) then $y^\otimes = \{y, y^*\}$. Then considering INLA notation it is obtained:

$$\begin{aligned}
\pi(y^* | y) &= \frac{\pi(y, y^*)}{\pi(y)} \text{ from the conditional probability} \\
&= \frac{\int \pi(y^* | \theta) \pi(y | \theta) \pi(\theta) d\theta}{\pi(y)} \text{ by exchangeability} \\
&= \frac{\int \pi(y^* | \theta) \pi(\theta | y) \pi(y) d\theta}{\pi(y)} \text{ applying Bayes' theorem} \\
&= \int \pi(y^* | \theta) \pi(\theta | y) d\theta
\end{aligned}$$

A DAG representation might offer the intuition behind Prediction in spatial models:

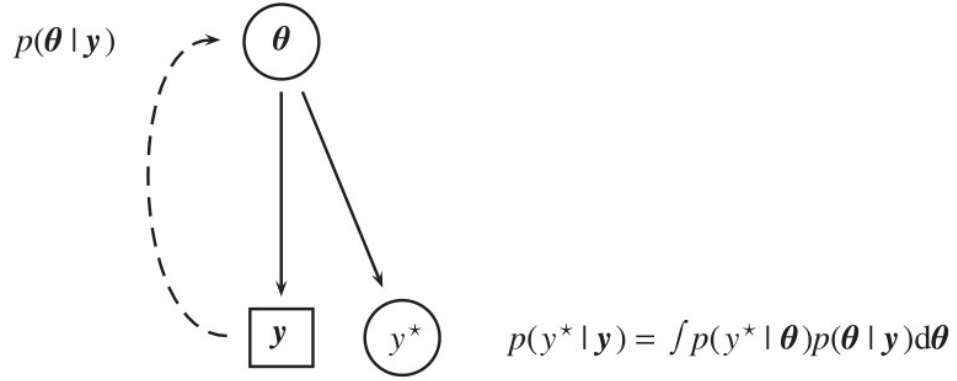


Figure 5.7: Spatial prediction representation through DAG, source Marta Blangiardo (2015)

where $\pi(y^* | y)$ is said predictive distribution and it is meaningful only in the Bayesian framework since the posterior distribution is treated as a random variable, which is totally not true in frequentist statistics.

5.8 Model Checking

(Incrociarlo con altri tesi)

Once the model is set up and fitted a resampling scheme has to be chosen in order to evaluate the model performance. One of the most used method to

assess bayesian model quality is LOOCV cross validation and default choice fo R-INLA package. From data is left out one single observation and so that the Validation set is $y_v = y_{-i}$ and the Assessement set is a $y_a = y_i$ the rest of the observations. Two KPI are assumed to be representative:

- CPO conditional predictive ordinate (pettit, 1990): $CPO_i = \pi(y^* | y_v)$
- PIT probability integral tranform (dawid, 1984): $PIT_i = \pi(y^* < y_i | y_v)$

These quantities are used by default by setting control options in the `inla(control.compute = list())` list object by setting them equal o TRUE. Inla also provides an inner method to authomatically handlee failing in computing those two quantities, leadind to values of 1 when predictions are not reliable and the ipposite for 0. Moreover the empirical distribution of the PIT can be used to asses predictive performance: if it is Uniform, so there are not values that strongly differ from the others then the model is correctly checked. Otherwise if the dostribtuon almost approxiamtes any of the other possibles then the Cross validation assesement prediction has led incorrectly predict the “out of the bag” validation sample.

Posterior checking method exploits a full cross validation where $y_a = y_v$ and it is called predictive checks. Th assesement set now is equal to the validation set, a s a consequence all the observation are evaluated twice. 4 quantities are driver to model estimate quality:

- the *posterior predictive distribution*: $\pi(y^* | y) = \int \pi(y^* | \theta_i) \pi(\theta_i | y) d\theta_i$ which is the likelihood of a replicate observation. When values are small that indicates that are those values are coming from tails, since the area under the curve (i.e. probability) is less. If this happens for many observation then outliers are driving the model leading to poor estimates
- the *posterior predictive p-value* whose math expression is: $\pi(y^* \leq y_i | y)$ for which values near to 0 and 1 indicates poor performances.
- *Root Mean Square Predictive Error RMSE*: $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - y_i^*)^2}$
- R^2

R-INLA has already anticipated in chapter 4 section 4.3.2 have designed a function to compute statistics on posterior distribution as `inla.pmarginal()` returning the cumulative density distribution.

5.9 Prior Specification

Chapter 6

SPDE approach

Observations in the spatial problem setting are considered as realizations of a stationary, isotropic unobserved GP $w(s)$ that we aim to estimate (5.1).

Before approaching the problem with SPDE, GPs were treated as multivariate Gaussian densities and Cholesky factorizations were applied on the covariance matrices and then fitted with likelihood. Matrices in this settings were very dense and they were scaling with the order of $O(n^3)$, leading to obvious big-n problem. The breakthrough, came with Lindgren et al. (2011) that proves that a stationary, isotropic GP with Matérn covariance can be represented as a GMRF using SPDE solutions by finite element method (Krainski, 2019). In other words given a GP whose covariance matrix is Q , SPDE can provide a method to approximate Q without computational constraints. As a matter of fact SPDE are equations whose solutions are GPs with a chosen covariance function focused on satisfying the relationship SPDE specifies. Benefits are many but the most important is that the representation of the GP through a GMRF provides a sparse representation of the spatial effect through a sparse precision matrix Q^{-1} . Sparse matrices enable convenient inner computation properties of GMRF that can be exploited with INLA. Bayesian inference on GMRF can take advantage of lower computational cost because of these properties stated before leading to a more feasible big-O $O(n^{3/2})$. The following chapter will

provide a intuition on SPDE oriented to practitioners. The chapter once again will follow the track of Krainski & Rubio (2019) and Blangiardo and Cameletti (2015) works, together with the street-opener paper from Miller et al. (2019) as compendium. SPDE might be complex for those who are not used to applied mathematics and physics making it difficult not only to grab the concept, but also to find its applications. One more obstacle regards SPDE software implementation, since without deep technical expertise it might be difficult to customize code with the aim to extend the methodology to different models. For a gentle introduction on what a SPDE is from a mathematical perspective a valuable reference is Miller et al. (2019) in section 2.1, then also its application to Matérn in 2.3.

6.1 Set SPDE Problem

Given the statistical model already encountered in chapter 5.4:

$$y(\mathbf{s}_i) = \mathbf{x}(\mathbf{s}_i)' \beta_j + w(\mathbf{s}) + \varepsilon(\mathbf{s}_i)$$

where $\eta(\mathbf{s}_i) = g(\mathbf{x}(\mathbf{s}_i)' \beta_j)$ is the linear predictor, whose link function $g(\cdot)$ is identity (can be also extended to GLM), where $w(\mathbf{s})$ is a Gaussian Process with mean structure 0 and $C(\cdot)$ covariance structure (where Q is the covariance matrix and Q^{-1} precision matrix). Then $w(s) \sim MVN(0, Q_{i,j}^{-1})$ and where $\varepsilon(\mathbf{s}_i)$ is white noise error such that $\varepsilon(\mathbf{s}_i) \sim \mathcal{N}(0, \tau^2)$.

Comprehending w in the model brings two major issues, specify a covariance function for observations as well as how to fit the model. Among all the possible reachable solutions including the SPDE, the common goal is to define covariance function between locations by approximating the precision matrix Q^{-1} , since they are an effective tool to represent covariance function as in section 4.1. For those reasons SPDE approach implies finding an SPDE whose solution have the precision matrix, that is desired for w . Lindgren

et al. (2011) proves that an approximate solution to SPDE equations is to represent w as a sum of basis function multiplied by coefficients (2019). Moreover the basis function coefficients are in reality a GMRF (for which fast method computations already exists).

6.2 SPDE within R-INLA

First point addresses the assumption that a GP with Matérn covariance function and $\nu > 0$ is a solution to *SPDE* equations. Second point addressed the issues of solving SPDE when grids are irregular, as opposite with the one seen in first point (regular grid for irregular distribution. In here comes FEM used in mathematics and engineering application with the purpose to solve differential equations. Notation is kept coherent with the one for the previous chapter.

6.3 First Point Krainsky Rubio TOO TECHNICAL

A regular 2D grid lattice is considered with infinite number of location points as vertices.

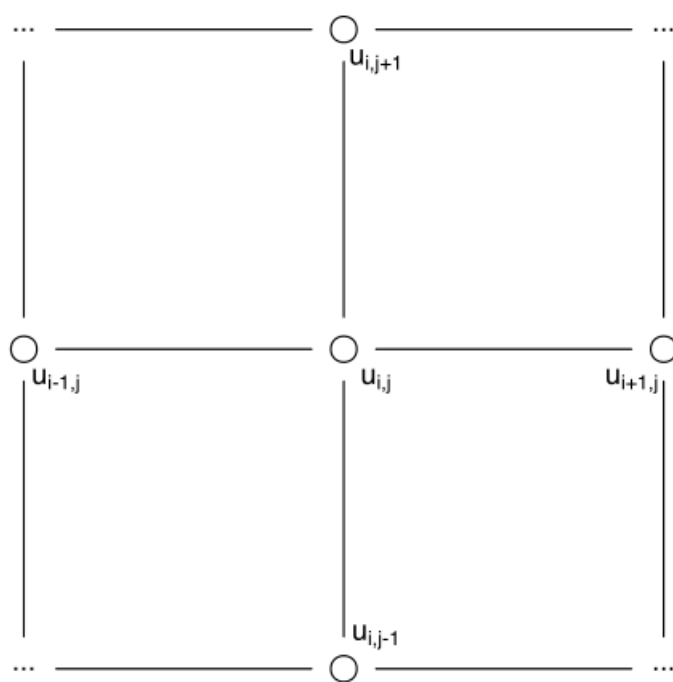


Figure 6.1: lattice 2D regular grid

Chapter 7

Exploratory Analysis

Data comes packed into the REST API end point `*/complete` in .json format. Data can be filtered out On the basis of the options set in the API endpoint argument body. Some of the options might regard the `city` in which it is evaluated the real estate market, `npages` as the number of pages to be scraped, `type` as the choice between rental of selling. For further documentation on how to structure the API endpoint query refer to section 3.3.3. Since to the analysis purposes data should come from the same source (e.g. Milan rental real estate within “circonvallazione”) a dedicated endpoint boolean option `.thesis` is passed in the argument body. What the API option under the hood is doing is specifying a structured and already filtered URL to be passed to the scraping endpoint. By securing the same URL to the scraping functions data is forced to come from the same URL source. The idea behind this concept can be thought as refreshing everyday the same immobiliare.it URL. API endpoint by default also specifies 10 pages to be scraped, in this case 120 is provided leading to 3000 data points. The *

refers to the EC2 public DNS that is

```
ec2-18-224-40-67.us-east-2.compute.amazonaws.com
```

```
http://*/complete/120/milano/affitto/.thesis=true
```

As a further source data can also be accessed through the mongoDB credentials with the cloud ATLAS database by picking up the latest .csv file

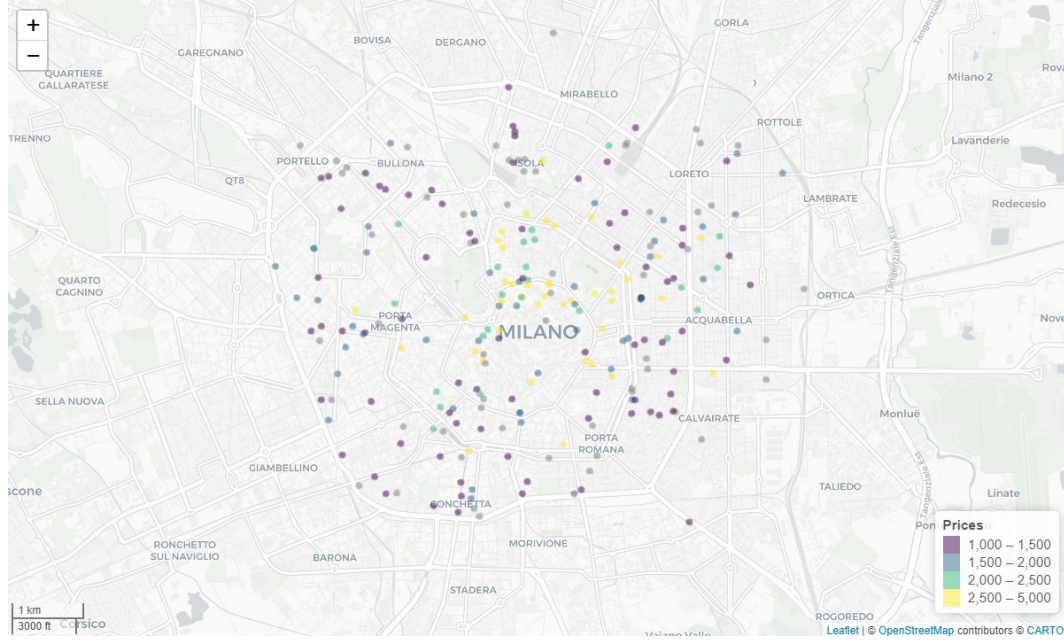
generated. For run time reasons also related to the bookdown files continuous building the API endpoint is called the day before the presentation so that the latest .csv file is available. As a consequence code chunks outputs are all cached due to heavy computation. Interactive maps are done with Leaflet, the result of which is a leaflet map object which can be piped to other leaflet functions. This permits multiple map layers and many control options to be added interactively (LaTeX output is statically generated)

A preliminary exploratory analysis evidences 34 covariates and 250 rows. Immobiliare.it furnishes many information regarding property attributes and estate agency circumstances. Data displays many NA in some of the columns but georeference coordinates, due to the design of scraping functions, are in any case present.

name	ref
ID	ID of the apartments
LAT	latitude coordinate
LONG	longitude coordinate
LOCATION	the complete address: street name and number
CONDOM	the condominium monthly expenses
BUILDAGE	the age in which the building was constructed
FLOOR	the floor the apartment is
INDIVSAPT	independent propeerty versus apartment
LOCALI	specification of the type and number of rooms
TPPROP	property type residential or not
STATUS	the actual status of the house, ristrutturato, nuovo, abitabile
HEATING	the heating system centralized or autonomous
AC	air conditioning hot and cold
PUB_DATE	the date of publication of the advertisement
CATASTINFO	land registry infromation
APTCHAR	apartment main characteristics

PHOTOSNUM	number of photos displays
AGE	estate agency name
LOWRDPRICE.originalPrice	If the price is lowered it indicates the starting price
LOWRDPRICE.currentPrice	If the price is lowered it indicates the current price
LOWRDPRICE.passedDays	If the price is lowered the days passed since the price has cha
LOWRDPRICE.date	If the price is lowered the date passed since the price has cha
ENCLASS	the energy class communicated to the
CONTR	the type of contract
DISP	if it is still avaiable or not
TOTPIANI	the total number of floors
PAUTO	number of parking box or garages avaiable in the property
REVIEW	estate agency review, long chr string
HASMULTI	it if has multimedia option, such as 3D house vitualization ho
PRICE	the monthly price <- response

Geographic coordinates can be represented on a map in order to get a first perception of spatial autocorrelations clusters. Leaflet maps are created with `leaflet()`, the result of which is a leaflet map object which can be piped to other leaflet functions. This allows multiple map layers and control settings to be added interactively. `leaflet` object takes as input data in latitude and longitude UTM coordinates so no transformation is required. Otherwise a projection to the right zone would be required and the `sp` transform



7.1 Data preparation

As already pointed out some data went missing since immobiliare provides data that in turn is filled by estate agencies or privates through pre compiled standard formats. Some of the missing observations can be reverse engineered by other information in the web pages e.g. given the street address it is possible to trace back the lat and long coordinates, even though this is already handled by the API through scraping functions with the mechanism found in section ???. Some of the information lacking in the summary table they might be desumed and then imputed by the estate agency review which is one of the covariates and where the most of the times missing information can be found. The approach followed in this part is to prune redundant information and “over” missing covariates trying to limit the dimensionality of the dataset.

7.1.1 NA Removal and Imputation

The first problem to assess is why information are missing. As already pointed out in the introduction part and in section ?? many of the presumably important covariates (i.e. price lat, long, title ,id ...) undergo to a sequence of forced step inside scraping functions with the aim to avoid to be missed. If in the end of the sequence the covariates values are still missing, the correspondent observation is not imputed and it is left out of the scraped dataset. The author choice to follow this approach relies on empirical observation that suggest when important information is missing then the rest of the covariates also do, as a consequence the observation is not useful. Moreover when the spatial component can no be reverse engineered then it is also the case of left out observation while scraping. The model needs to have a spatial component in order to be evaluated. To this purpose The missing profile is crucial since it can also suggest problems in the scraping procedure.

By identifying pattern in missing observation the maintainer can take advantage of it and then debug the part that causes missingness. In order to identify the pattern a revision of missingness is introduced by *Little and Rubin (2014)* *miss lit* .randomnes can be divided into 3 categories:

- MCAR (missing completely at random) likelihood of missing is equal for all the information, in other words missing data are one independent for the other.
- MAR (missing at random) likelihood of missing is not equal.
- *NMAR* (not missing at random) data that is missing due to a specific cause, scraping can be the cause.

The last hypothesis MNAR is often the case of daily monitoring clinical studies, where patient might drop out the experiment because of death and so as a consequence all the data starting from the death time +1 is missing.

To assess pattern A *heat map* plot fits the need:

Looking at the left hand of the heat map plot 7.1 considering from

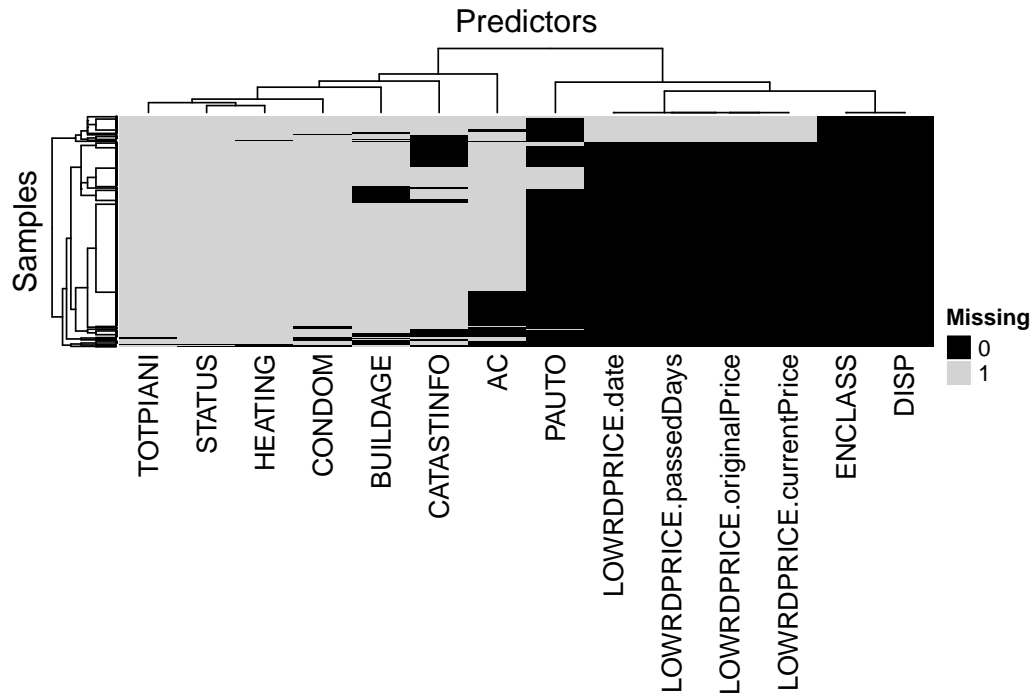


Figure 7.1: Missingness Heatmap plot

TOTPIANI to *AC* there are no relevant patterns and missingness can be traced back to MAR, conditioned mean imputation is applied until *CONDOM* included, the others are discarded. In the far right hand side *ENCLASS* and *DISP* data is completely missing, this can be addressed to a scraping fail. Further inspection of the API scraping process focused on those covariates is strongly advised. From *LOWRDPRICE.* covariates class it seems to be witnessing a missing underlining pattern NMAR which is clearer by looking at the second co_occurrence plot 7.2 analysis. Co-occurrence analysis might suggest frequency of missing predictor combinations and *LOWRDPRICE.* class covariates are missing all together in combination. *PAUTO* is missing where lowered price class covariates are missing but this is not true for all the observations leading to the conclusion that *PAUTO* should be treated as a low prevalence covariate, therefore *PAUTO* is discarded. After a further investigation *LOWRDPRICE.* exists when the price covariates is effectively decreased, that leads to group the covariate's

information and to encode it as a two levels categorical covariate. Further methods to feature engineer the lowrdprice covariates are profile data.

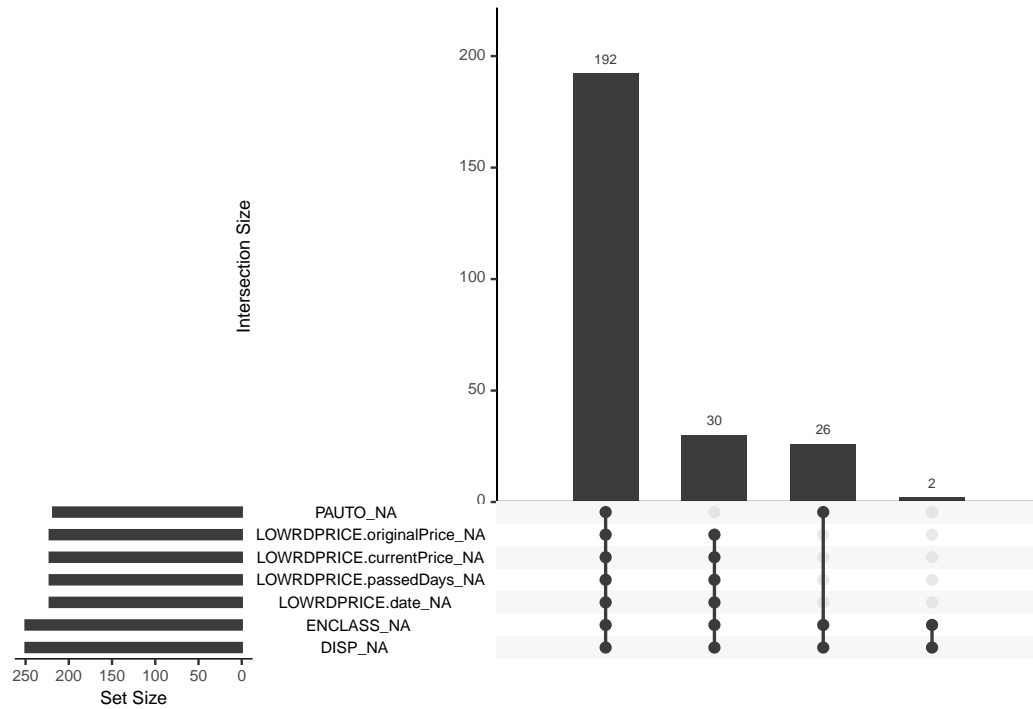


Figure 7.2: Missingness co-occurrence plot

7.2 Spatial Autocorrelation assesement

Spatial data come packed into point reference

- tmap
- leaflet
- gganimate

7.3 Model Specification

7.4 Mesh building

PARAFRASARE The SPDE approach approximates the continuous Gaussian field w_i as a discrete Gaussian Markov random field by means of a finite basis function defined on a triangulated mesh of the region of study. The spatial surface can be interpolated performing this approximation with the `inla.mesh.2d()` function of the R-INLA package. This function creates a Constrained Refined Delaunay Triangulation (CRDT) over the study region, that will be simply referred to as the mesh. Mesh should be intended as a trade off between the accuracy of the GMRF surface representation and the computational cost, in other words the more are the vertices, the finer is the GF approximation, leading to a computational funnel.

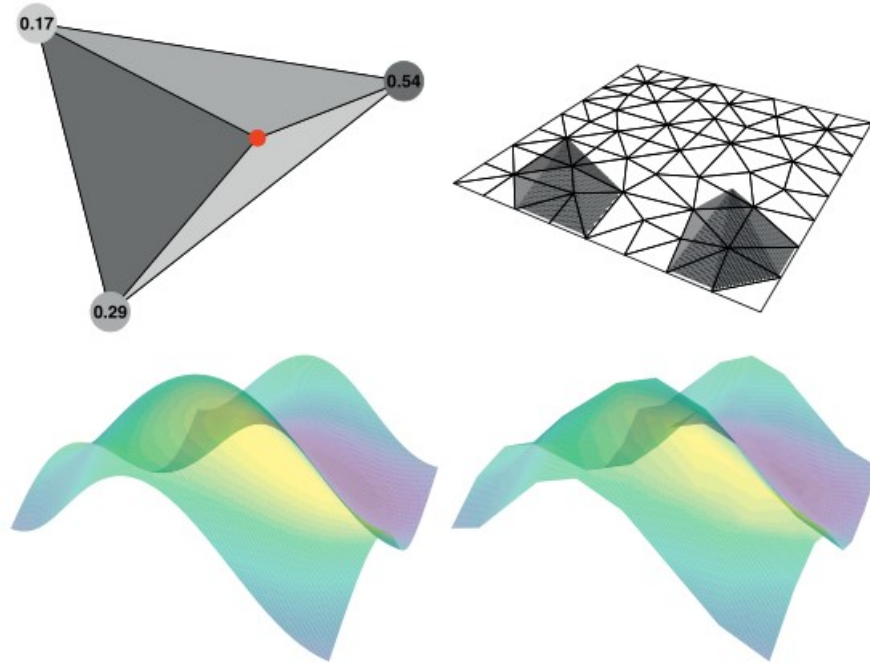


Figure 7.3: Triangularization intuition, Krainski (2019) source

Arguments can tune triangularization through `inla.mesh.2d()` :

- `loc`: location coordinates that are used as initial mesh vertices

- `boundary`: object describing the boundary of the domain,
- `offset`: argument is a numeric value (or a length two vector) and it is used to set the automatic extension distance. If positive, it is the extension distance in the same scale units. If negative, it is interpreted as a factor relative to the approximate data diameter; i.e., a value of -0.10 (the default) will add a 10% of the data diameter as outer extension.
- `cutoff`: points at a closer distance than the supplied value are replaced by a single vertex. Hence, it avoids small triangles
- `max.edge`: A good mesh needs to have triangles as regular as possible in size and shape.
- `min.angle` argument (which can be scalar or length two vector) can be used to specify the minimum internal angles of the triangles in the inner domain and the outer extension

A convex hull is a polygon of triangles out of the domain area, in other words the extension made to avoid the boundary effect. All meshes in Figure 2.12 have been made to have a convex hull boundary. If borders are available are generally preferred, so non convex hull meshes are avoided.

7.4.1 Shinyapp for mesh assessment

INLA includes a Shiny (Chang et al., 2018) application that can be used to tune the mesh params interactively

The mesh builder has a number of options to define the mesh on the left side.

These include options to be passed to functions `inla.nonconvex.hull()` and `inla.mesh.2d()` and the resulting mesh displayed on the right part.

7.4.2 BUilding SPDE model on mesh

7.5 Spatial Kriging (Prediction)

QUI INCERTEZZE

Chapter 8

Model Selection & Fitting

8.1 Model Criticism

evaluation of the variables to include in the model,, assumptions of the model
i.e. exchangeability and independence prior distribution to assign to
parameters and hyper parameters.

8.2 Spatial Kriging

8.3 Model Checking

if (models > 2){

Model Selection

IDEA: proporre due modelli uno più interpretabile con distribuzione normale

}

Chapter 9

Final Words

We have finished a nice book.

Bibliography

(2018).

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.

Banerjee, S., Carlin, B. P., and Gelfand, A. E. (2014). *Hierarchical Modeling and Analysis for Spatial Data*. Chapman and Hall/CRC.

Bell, B. S., Hoskins, R. E., Pickle, L., and Wartenberg, D. (2006). *International Journal of Health Geographics*, 5(1):49.

Blanchet-Scalliet, C., Helbert, C., Ribaud, M., and Vial, C. (2019). Four algorithms to construct a sparse kriging kernel for dimensionality reduction. *Computational Statistics*, pages 1–21.

Cameletti, M., Lindgren, F., Simpson, D., and Rue, H. (2012). Spatio-temporal modeling of particulate matter concentration through the SPDE approach. *ASTA Advances in Statistical Analysis*, 97(2):109–131.

Capozza, D. and Seguin, P. (1996). Expectations, efficiency, and euphoria in the housing market. *Regional Science and Urban Economics*, 26:369–386.

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.

Clark, T. E. (1995). Rents and prices of housing across areas of the United States. A cross-section examination of the present value model. *Regional Science and Urban Economics*, 25(2):237–247.

- Cressie, N. (2015). *Statistics for Spatial Data*. Probability and Mathematical Statistics. Wiley-Interscience, revised edition edition.
- Densmore, J. (2019). Ethics in web scraping.
- Diestel, R. (2006). *Graph theory*. Graduate Texts in Mathematics. Springer, 3rd edition.
- Google (2020). Presentazione dei file robots.txt - guida di search console.
- Gómez Rubio, V. (2020). *Bayesian Inference with INLA*. Chapman and Hall/CRC.
- Herath, S. and Maier, G. (2011). Hedonic house prices in the presence of spatial and temporal dynamics. *Territorio Italia - Land Administration Cadastre, Real Estate*, 1:39–49.
- Inc., D. (2020a). Get docker.
- Inc., R. H. (2020b). 7.2. advantages of using docker red hat enterprise linux 7.
- Khalil, S. (2018). *Rcrawler: Web Crawler and Scraper*. R package version 0.1.9-1.
- Krainski, E., Gómez-Rubio, V., Bakka, H., Lenzi, A., Castro-Camilo, D., Simpson, D., Lindgren, F., and Rue, H. (2018). *Advanced Spatial Modeling with Stochastic Partial Differential Equations Using R and INLA*. Chapman and Hall/CRC.
- Krainski, E. T. (2019). *Advanced spatial modeling with stochastic partial differential equations using R and INLA*. Chapman and Hall/CRC.
- Lancaster, K. J. (1966). A new approach to consumer theory. *Journal of Political Economy*, 74(2):132–157.
- Li, H., Li, H., Ding, Z., Hu, Z., Chen, F., Wang, K., Peng, Z., and Shen, H. (2020). Spatial statistical analysis of coronavirus disease 2019 (covid-19) in china. *Geospatial Health*, 15(1).

- Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between gaussian fields and gaussian markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(4):423–498.
- Ling, Y. and Ling, Y. (2019). Time, space and hedonic prediction accuracy evidence from the corsican apartment market. *The Annals of Regional Science*, page 28.
- Malpezzi, S. (2008). *Hedonic Pricing Models: A Selective and Applied Review*, pages 67–89.
- Manganelli, B., Morano, P., and Tajani, F. (2013). Economic relationships between selling and rental prices in the italian housing market.
- Marta Blangiardo, M. C. (2015). *Spatial and Spatio-temporal Bayesian Models with R-INLA*. Wiley.
- Meissner, P. and Ren, K. (2020). *robotstxt: A 'robots.txt' Parser and 'Web-bot'/'Spider'/'Crawler' Permissions Checker*. R package version 0.7.7.
- Microsoft and Weston, S. (2020). *foreach: Provides Foreach Looping Construct*. R package version 1.5.0.
- Miller, D. L., Glennie, R., and Seaton, A. E. (2019). Understanding the stochastic partial differential equation approach to smoothing. *Journal of Agricultural, Biological and Environmental Statistics*, 25(1):1–16.
- Moraga, P. (2019). *Geospatial Health Data*. Chapman and Hall/CRC.
- Nolis, J. (2020). R docker faster.
- Paci, L., Beamonte, M. A., Gelfand, A. E., Gargallo, P., and Salvador, M. (2017). Analysis of residential property sales using space–time point patterns. *Spatial Statistics*, 21:149 – 165.
- Perepolkin, D. (2019). *polite: Be Nice on the Web*. R package version 0.1.1.

- Rosen, S. (1974). Hedonic prices and implicit markets: Product differentiation in pure competition. *Journal of Political Economy*, 82(1):34–55.
- Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields*. Chapman and Hall/CRC.
- Rue, H., Martino, S., and Chopin, N. (2009). Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(2):319–392.
- Shah, T. (2018). *ratelimitr: Rate Limiting for R*. R package version 0.4.1.
- Trestle Technology, LLC (2018). *plumber: An API Generator for R*. R package version 0.4.6.
- Vaughan, D. and Dancho, M. (2018). *furrr: Apply Mapping Functions in Parallel using Futures*. R package version 0.1.0.
- was previously an Economist at the Indeed Hiring Lab, A. F. F. (2020). Indeed tech skills explorer: Today’s top tech skills.
- WhoIsHostingThis.com (2020). User agent: Learn your web browsers user agent now.
- Wickham, H. (2019). *rvest: Easily Harvest (Scrape) Web Pages*. R package version 0.3.5.
- Wickham, H., Hester, J., Müller, K., and Cook, D. (2017). *memoise: Memoisation of Functions*. R package version 1.1.0.
- 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.