

# Modulo 1 Geostatistica - corso base

Francesco Pirotti

2024-09-05

## Contents

<b>1</b>	<b>Introduzione ad R</b>	<b>3</b>
1.1	Cosa impareremo . . . . .	3
1.2	Materiale . . . . .	3
1.3	R: panoramica . . . . .	3
1.4	R: panoramica . . . . .	3
1.5	RStudio . . . . .	4
1.6	Righe codice - i comandi . . . . .	5
1.7	Commenti e sezioni . . . . .	5
1.8	Documentazione/help . . . . .	5
1.9	Packages/Librerie . . . . .	6
1.10	Packages/Librerie . . . . .	6
1.11	Variabili e funzioni . . . . .	6
1.12	Variabili e funzioni: Scope . . . . .	7
1.13	Strutture dati in R . . . . .	7
1.14	Tipo dati: vettori . . . . .	8
1.15	Tipo dati: matrix . . . . .	8
1.16	Tipo dati: array . . . . .	9
1.17	Tipo dati: List . . . . .	9
1.18	Tipo dati: List c/nomi . . . . .	10
1.19	Tipo dati: Data frame . . . . .	11
1.20	Tipo dati: Data frame . . . . .	11
1.21	Tipo dati: Tibble . . . . .	12
1.22	Convertire oggetti . . . . .	12
1.23	Salvare oggetti R . . . . .	13
1.24	Operatori di R . . . . .	13
1.25	Esercizio 1 . . . . .	14

<b>2</b>	<b>Data in/process/out</b>	<b>14</b>
2.1	Cosa faremo . . . . .	14
2.2	Input dati: CSV . . . . .	14
2.3	Input dati: Fogli excel . . . . .	15
2.4	Grafici di base . . . . .	15
2.5	Grafici di base . . . . .	16
2.6	Input dati: JSON . . . . .	17
<b>3</b>	<b>Trasformazione dei dati</b>	<b>17</b>
3.1	Introduzione . . . . .	17
3.2	Iteratori . . . . .	18
<b>4</b>	<b>Tidyverse</b>	<b>18</b>
4.1	Introduzione . . . . .	18
4.2	Iteratori: map . . . . .	18
4.3	Iteratori: map . . . . .	19
4.4	Iteratori: map . . . . .	19
4.5	Filter e select . . . . .	19
4.6	Mutate . . . . .	20
4.7	Aggregazioni . . . . .	20
4.8	Unione tabelle (join) . . . . .	20
4.9	Unione tabelle (join) . . . . .	21
4.10	Esercizio 2 . . . . .	22
4.11	Ristrutturazione: wide=>long . . . . .	22
4.12	Ristrutturazione: long=>wide . . . . .	23
<b>5</b>	<b>Grafici</b>	<b>23</b>
5.1	Grafici con GGplot - dispersione . . . . .	23
5.2	Grafici con GGplot - scatola e baffi . . . . .	24
5.3	Grafici multipli (facets) . . . . .	24
5.4	Grafici multipli (facets) . . . . .	25
5.5	Grafici e dati geo . . . . .	28
5.6	Grafici e stampe PNG . . . . .	29
5.7	Grafici e stampe PDF . . . . .	29
5.8	Cheatsheets . . . . .	30
5.9	Esercizio 3 . . . . .	30
<b>6</b>	<b>Dati geospaziali</b>	<b>31</b>
6.1	Introduzione . . . . .	31
6.2	Dati vettoriali . . . . .	31
6.3	Raster . . . . .	31
6.4	Interpolazione . . . . .	33

# 1 Introduzione ad R

## 1.1 Cosa impareremo

- struttura di R (base e pacchetti), potenzialità
- RStudio: come funziona (panoramica)
- documentazione, esempi, snippets per imparare
- principali categorie di dati: oggetti, variabili e funzioni in R
- assegnazione ed utilizzo variabili
- leggere, rielaborare e scrivere dati strutturati (tabelle) e non
- leggere ed interpretare dati *geo-spaziali*
- rappresentare dati con grafici e mappe

## 1.2 Materiale

Trovate il materiale in una pagina GitHub [QUI](#) e potrete scaricarlo anche come zip ed anche direttamente nello spazio web [QUI](#). Troviamo:

- slides in formato HTML
- contenuto slides in formato di documento A4 PDF
- materiale per gli esercizi

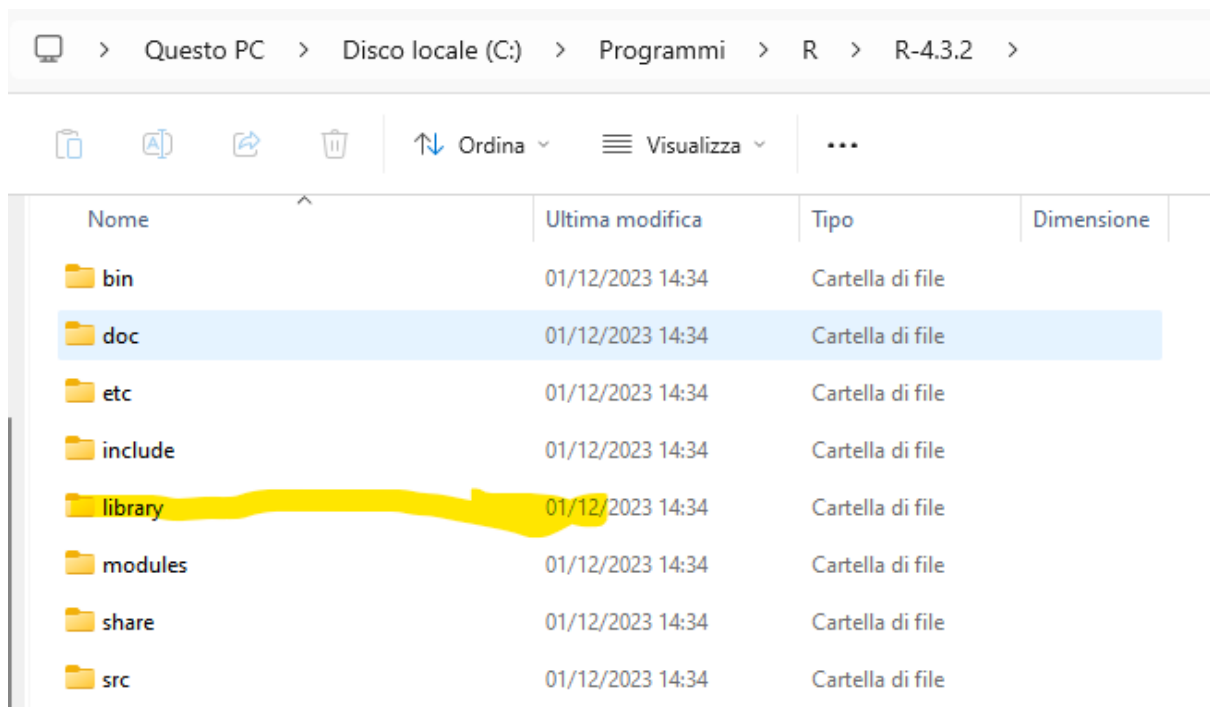
Verrà fornito un link GDrive per caricare gli esercizi per le verifiche.

## 1.3 R: panoramica

- Ma R serve solo per analisi statistiche? No!
  - E.g. questo materiale è creato con R (vedi *rmarkdown*)
  - Le analisi possono diventare interattive con *Shiny* (esempi progetti InforSAT, VARCITIES...)
  - ....

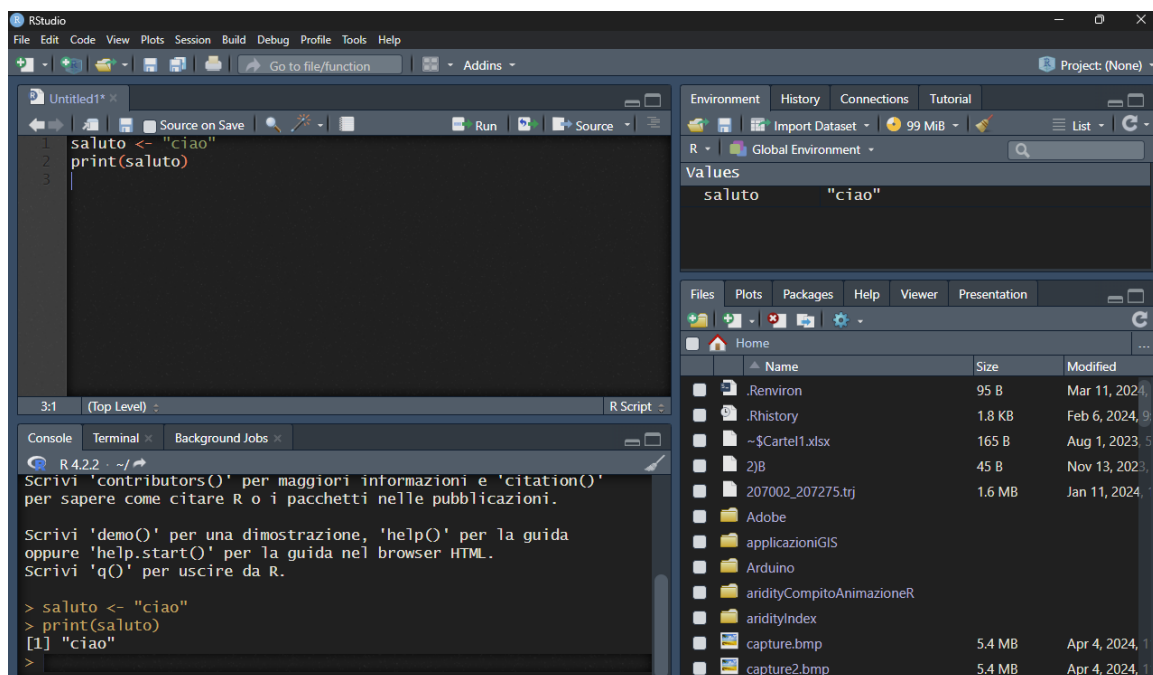
## 1.4 R: panoramica

- Installazione di R, versioni
- struttura di R - percorso di installazione



## 1.5 RStudio

- versione desktop/server
- vantaggi interfaccia:
  - console/terminale/background
  - progetti/packages/help
  - Environment/History/Connection/Tutorial
- Salvare un progetto, cartella di progetto, GIT e GitHub
- Working Directory predefinita, uso del tab per richiamare percorsi



## 1.6 Righe codice - i comandi

R è un programma basato su righe di codice con comandi che eseguono elaborazioni su dati.

L'utente immette i comandi al prompt ( `>` ) e ciascun comando viene eseguito uno alla volta andando a capo.

Le righe di comando solitamente vengono salvate in un file "script" con estensione "R" (.R) e vengono eseguite una alla volta mediante "invio" o con selezione multipla e "invio".

Con RStudio è possibile eseguire l'intero file, fermandosi eventualmente in punti specifici "breakpoints" (lo vedremo durante il corso).

Esercizio: esegui comando della figura alla slide precedente.

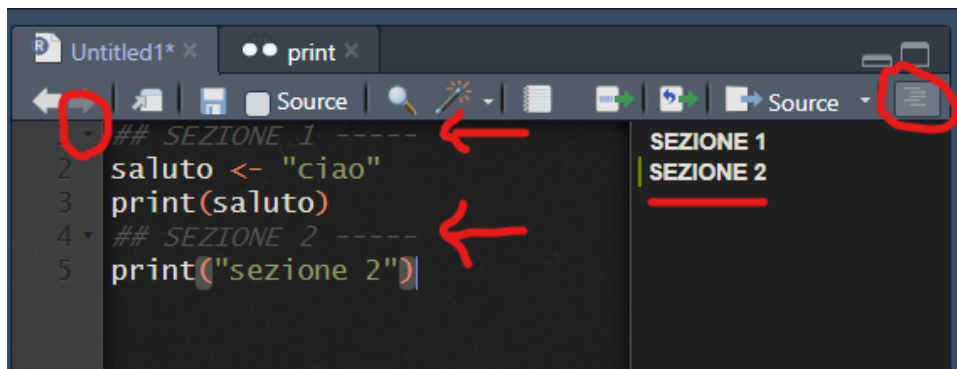
```
saluto <- "ciao"  
print(saluto)
```

```
## [1] "ciao"
```

## 1.7 Commenti e sezioni

Se non si vuole eseguire delle righe, basta mettere un carattere asterisco (#) all'inizio del testo che NON si vuole eseguire.

NB se inserito all'inizio della riga, e seguito da testo e da 4 o più caratteri - o #, il testo diventa una sezione



### 1.7.1

## 1.8 Documentazione/help

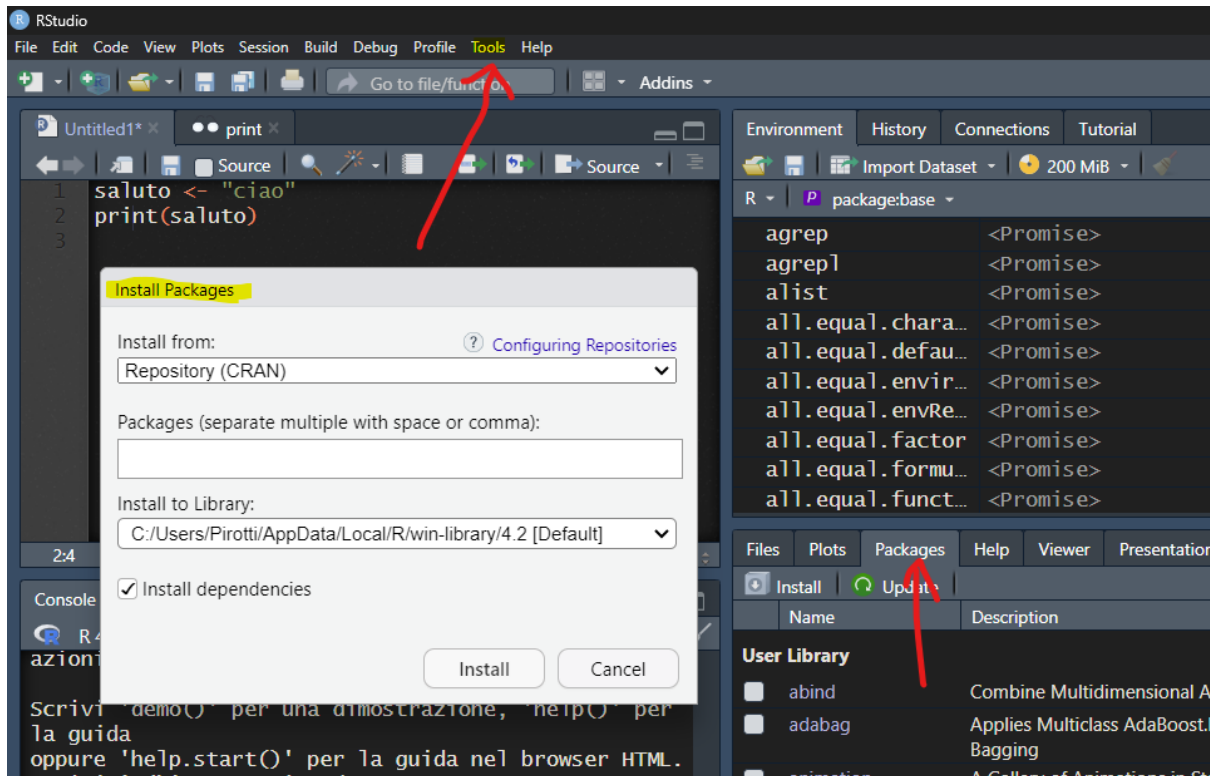
Ogni singola funzione ha ampia documentazione con molti esempi. Chiamando una funzione dopo uno o due punti interrogativi richiama la documentazione. Quasi sempre gli esempi sono eseguibili facendo copia/incolla

Esercizio: esegui il primo esempio dalla documentazione della funzione *print*

```
?print  
??print
```

## 1.9 Packages/Librerie

Tantissime funzionalità aggiuntive sono disponibili su componenti aggiuntivi che vanno installate con il comando `install.packages(<nome libreria>)` e poi caricate con il comando `library(<nome libreria>)`. Da RStudio possiamo caricarle da interfaccia grafica.



## 1.10 Packages/Librerie

Non tutti sono nella “repository” CRAN, alcuni sono su GitHub o scaricabili da altre fonti - vedi alcune cose sotto

- <https://github.com/fpirotti/rPET>
- <https://github.com/fpirotti/CloudGeometry>

Per installare librerie presenti solo in CRAN si utilizza una libreria chiamata “devtools”

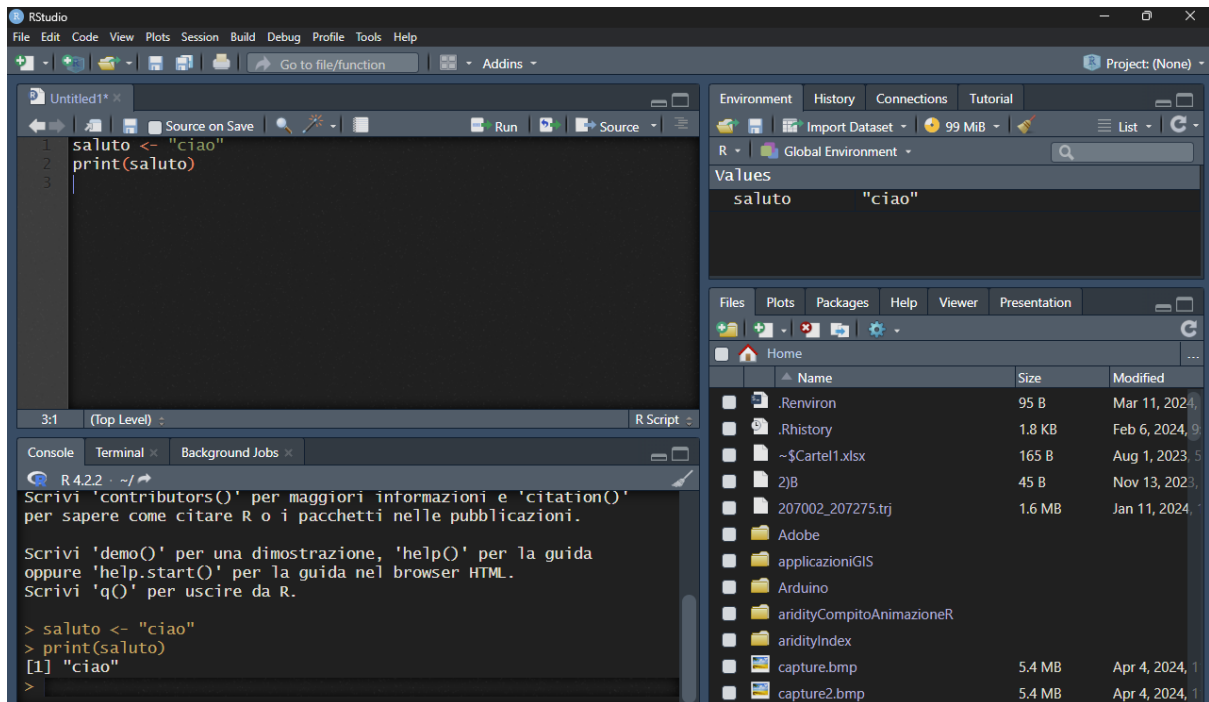
```
# install.packages("devtools")
devtools::install_github("fpirotti/CloudGeometry")
```

## 1.11 Variabili e funzioni

Qui vediamo una variabile ed una funzione.

**NB1** - la variabile “saluto” è nel Environment (“ambito/ambiente/campo”) *globale\**. La funzione “print” è nel campo del package “base” - tieni premuto il tasto CTRL e seleziona il nome della funzione - vedi cosa succede.

**NB2** - operatore di assegnazione <- (o <<- nel caso si voglia forzare l’assegnazione ad una variabile *globale\**)



## 1.12 Variabili e funzioni: Scope

\*Le variabili create al di fuori di funzioni sono note come variabili *globali*; possono essere utilizzate sia all'interno delle funzioni che all'esterno.

Sotto andiamo a creare una nostra funzione “salutami” che esegue il saluto. Provate a modificare l'operatore di assegnazione da <- a <<- e rieseguire!!

```
salutami <- function(){
  saluto <- "ciao ARPA!!!"
  print(saluto)
}
```

```
print(saluto)
```

```
## [1] "ciao"
```

```
salutami()
```

```
## [1] "ciao ARPA!!!"
```

```
print(saluto)
```

```
## [1] "ciao"
```

## 1.13 Strutture dati in R

NB ogni elemento in R è considerato (ed è) un VETTORE. Le funzioni di R considerano ogni variabile un vettore. Cosa significa? Che le funzioni elaborano tutti gli elementi di un vettore “by default” e che ogni elemento è indicabile con un numero iniziando da 1 (non da 0 come solitamente succede in altri linguaggi).

- vector
- character
- integer
- numeric

```
miaVar <- FALSE
class(miaVar)
miaVar[[1]]
miaVar[[2]]
```

## 1.14 Tipo dati: vettori

Perchè succede quello che vedete sotto (dati *numeric* diventano *character*)?

```
miaVar <- c(1,4,6,8)
class(miaVar)
```

```
## [1] "numeric"
```

```
miaVar[[1]]
```

```
## [1] 1
```

```
miaVar[[2]]
```

```
## [1] 4
```

```
miaVar[[2]] <- "evviva"
class(miaVar)
```

```
## [1] "character"
```

```
print(miaVar)
```

```
## [1] "1"      "evviva" "6"      "8"
```

## 1.15 Tipo dati: matrix

Matrix è un oggetto con struttura di matrice ovvero bidimensionale (righe  $\times$  colonne); pensate a un gruppo di vettori impilati o affiancati.

Si accede e si assegnano i valori con  $[r,c]$  dove  $r$  e  $c$  sono gli indici di riga e colonna. Si può lasciare vuoto un indice per accedere alla riga/colonna

```
mat <- matrix(c(1, 2, 3, 4), nrow = 2, ncol = 2)
mat[[2]]
```

```
## [1] 2
```

```
mat[1,2]
```

```
## [1] 3
```



```
mat[1,]
```

```
## [1] 1 3
```

```
mat[1,2] <- 100
```

## 1.16 Tipo dati: array

Un array è una matrice multidimensionale.

r=rows, c=columns, m=matrice... etc...

Esercizio: vedi sotto come creare un array a 3 dimensioni. Nota che duplica 9 valori 2 volte. Prova a dare 8 valori invece che nove. Prova a dare 3 valori. Cosa succede.

```
# 2 vettori di valori
valori1 <- c(5, 9, 3)
valori2 <- c(10, 11, 12, 13, 14, 15)
column.names <- c("C1", "C2", "C3")
row.names <- c("R1", "R2", "R3")
matrix.names <- c("Matrix1", "Matrix2")

# Crea un array a tre dimensioni
arr <- array(c(valori1, valori2), dim = c(3, 3, 2),
             dimnames = list(row.names,
                              column.names,
                              matrix.names))

print(arr)
```

```
## , , Matrix1
##
##      C1 C2 C3
## R1   5 10 13
## R2   9 11 14
## R3   3 12 15
##
## , , Matrix2
##
##      C1 C2 C3
## R1   5 10 13
## R2   9 11 14
## R3   3 12 15
```

## 1.17 Tipo dati: List

Le strutture vector/matrix/array, possono contenere solo una tipologia base (numeric, integer, character, logical...). Ma la struttura LIST no!

La struttura *list* è un set di dati eterogenei; opzionalmente, è possibile assegnare dei nomi a ciascun elemento nel set.

```
lista <- list(1,4,6,8)
class(lista)
```

```
## [1] "list"
```

```
lista[[1]]
```

```
## [1] 1
```

```
lista[[2]]
```

```
## [1] 4
```

```
lista[[2]] <- "evviva"  
class(lista)
```

```
## [1] "list"
```

```
lista[[1]]
```

```
## [1] 1
```

```
lista[[2]]
```

```
## [1] "evviva"
```

```
print(lista)
```

```
## [[1]]  
## [1] 1  
##  
## [[2]]  
## [1] "evviva"  
##  
## [[3]]  
## [1] 6  
##  
## [[4]]  
## [1] 8
```

## 1.18 Tipo dati: List c/nomi

NB, la struttura LIST non è altro che un set, una “lista”, di oggetti associata ad un indice. L’indice è un numero intero ma può essere un testo (simile al concetto di coppie “key->value”/chiave->valore).

Può essere assegnato un indice/chiave per riferimento all’elemento nella lista

```
names(lista)
```

```
## NULL
```

```
names(lista) <- c("primo", "secondo")  
lista[[1]]
```

```
## [1] 1
```

```
lista[["primo"]]
```

```
## [1] 1
```

Svantaggi di strutture tipo *list*: usa più memoria!

**Nota su allocazione memoria:** se dovete gestire volumi importanti di dati, considerate la pre-allocazione della memoria SE conoscete la dimensione. Vedi anche blog qui.

```
vettoreMoltoGrande = numeric(1000)
vettoreMoltoGrande[[100]]
```

```
## [1] 0
```

## 1.19 Tipo dati: Data frame

Un *data frame* è una struttura di tipo *list* ma con un numero uguale di “righe” per ogni colonna di attributi. È possibile manipolare i data frame filtrando sulle righe e operando sulle colonne.

Sia righe che colonne possono avere identificativi.

```
c.lat <- c(45.1, 45.2, 45.3)
c.lon <- c(11.1, 11.2, 11.3)

df.geo.stz <- data.frame(stz=c("A", "B", "C"),
                        longitude=c.lon,
                        latitudine=c.lat)

colnames(df.geo.stz)
```

```
## [1] "stz"          "longitude" "latitudine"
```

```
rownames(df.geo.stz)
```

```
## [1] "1" "2" "3"
```

```
rownames(df.geo.stz) <- df.geo.stz$stz
rownames(df.geo.stz)
```

```
## [1] "A" "B" "C"
```

## 1.20 Tipo dati: Data frame

Le colonne sono *vettori* - si può richiamare e assegnare i valori di una colonna con `$` o `[[ ]]` o `[, "<nomecolonna>"]` o `[ riga, "<nomecolonna>"]`

Come con Tipo dati: matrix si possono usare gli indici di riga/colonna.

```
df.geo.stz$stz
df.geo.stz[, "stz"]
df.geo.stz[["stz"]]
df.geo.stz$stz <- c("A1", "B1", "C1")
```

## 1.21 Tipo dati: Tibble

Un data frame particolare, lo vedremo quando usiamo l'infrastruttura di librerie “tidyverse”. Vantaggi:

- Non modifica il nome mettendo un punto al posto degli spazi (Vedi esempio)
- Non utilizza mai `row.names()`. Lo scopo di *tidyverse* è quello di memorizzare le variabili in modo coerente. Quindi non memorizza mai una variabile come attributo speciale (i.e. nomi riga).

```
library(tidyverse)
geo.stz.tb <- tibble(df.geo.stz)
class(geo.stz.tb)
```

```
## [1] "tbl_df"      "tbl"        "data.frame"
```

```
#geo.stz.tb
names(data.frame(`nome colonna con spazi` = 1))
```

```
## [1] "nome.colonna.con.spazi"
```

```
names(tibble(`nome colonna con spazi` = 1))
```

```
## [1] "nome colonna con spazi"
```

- Ricicla solo vettori di lunghezza 1. Questo perché il riciclo di vettori di lunghezza superiore è una fonte frequente di problemi nascosti. Riciclo: prova i seguenti comandi:

```
v <- 1:10
v[]
```

```
## [1] 1 2 3 4 5 6 7 8 9 10
```

```
v[] <- 1
v[]
```

```
## [1] 1 1 1 1 1 1 1 1 1 1
```

```
v[] <- 1:5
v[]
```

```
## [1] 1 2 3 4 5 1 2 3 4 5
```

```
# sotto non funziona!
#v[] <- 1:3
```

## 1.22 Convertire oggetti

Oggetti tra loro compatibili si possono convertire; ad esempio tra `matrix` => `data.frame` usando il comando `as.data.frame(<oggetto matrix>)` viceversa attenzione che tutti i tipi di vettore vengono resi omogenei.

Provate il seguente esercizio sotto.

```
mat
```

```
##      [,1] [,2]
## [1,]    1 100
## [2,]    2   4
```

```
df.mat <- as.data.frame(mat)
as.matrix(df.mat)
```

```
##      V1 V2
## [1,]  1 100
## [2,]  2   4
```

```
df.mat$nuovaColonna <- "testo"
as.matrix(df.mat)
```

```
##      V1 V2  nuovaColonna
## [1,] "1" "100" "testo"
## [2,] "2" " 4" "testo"
```

```
df.mat
```

```
##   V1 V2 nuovaColonna
## 1  1 100          testo
## 2  2  4          testo
```

## 1.23 Salvare oggetti R

Rstudio può salvare l'intero progetto con le variabili e le funzioni che vedete in alto a destra.

Il comando `save` salva in un file con estensione “rda” che viene riconosciuto anche direttamente da RStudio. NB è possibile salvare qualsiasi oggetto, anche funzioni!

Prova a cliccare sul file dopo aver lanciato la prima riga del comando seguente!

```
## NB salvo anche la funzione "salutami" fatta all'inizio!
save(df.geo.stz, salutami, file="oggetti.rda")
load("oggetti.rda")
saveRDS(df.geo.stz, file = "geo.stz.RDS")
df.geo.stz <- readRDS("geo.stz.RDS")
```

## 1.24 Operatori di R

Gli operatori aritmetici e logici di R funzionano sia su singoli scalari che su vettori e strutture come *array* e *matrix*.

NB questo vuol dire che si possono eseguire operazioni su tutti i singoli valori della struttura internamente (*vectorization*)!

```
c.lat + 10
c.lat == c.lon

c.lat[[1]] <- 44
c.lat <- c.lat - 1
## AND / OR
```

```
c.lat == c.lon & c.lat==0
c.lat == c.lon | c.lat==45.1
df.mat$V3 <- df.mat$V2**2
#c.lat == c.lon && c.lat==c.lat
#c.lat == c.lon || c.lat==c.lat
```

## 1.25 Esercizio 1

- Creare una funzione chiamata *crea.df* che esegue la seguente operazione:
  - creazione di un oggetto data.frame con 2 colonne, (Stazione, ValoreX) e 4 righe. Inventate voi i valori. Calcolo una terza colonna con il quadrato di ValoreX.
- Eseguire la funzione per creare un oggetto di tipo “tibble” chiamato *tabella.tb*
- **Consegnate** un file RDA contenente i due oggetti, la funzione *crea.df* e il tibble *tabella.tb* - chiamate il file con: *es1\_cognome\_nome.rda*, e.g. *es1\_pirotti\_francesco.rda*

## 2 Data in/process/out

### 2.1 Cosa faremo

Impareremo a leggere/trasformare/scrivere dati da fonti diverse:

- file di testo strutturati (**CSV**)
- fogli di lavoro (**XLSX**)
- dati online non strutturati come **JSON**

### 2.2 Input dati: CSV

I file CSV sono notoriamente file di testo strutturati come tabelle e molto utilizzati per condividere dati.

Il comando standard è `read.csv`.

NB non è ottimizzato per l'efficienza. Per dati massivi usare il comando `fread` dalla libreria `data.table` `data.table::fread()`.

Facciamo un esercizio con dati ARPA! Scarichiamo dalla pagina [QUI](#) o direttamente da [QUI](#) i dati di precipitazione mensile - salvateli in una sottocartella “dati” e caricateli usando sia la funzione `read.csv` che `fread` dalla libreria `data.table` (se l'avete altrimenti vedete a schermo i risultati). Si vuole notare le differenze e come si chiama una funzione SENZA caricare una libreria, usando `::`.

```
prec <- read.csv("dati/Prec_mensili_2020.csv")
prec2 <- data.table::fread("dati/Prec_mensili_2020.csv")
summary(prec2)
```

```
## Precipitazioni mensili Veneto 2020 (mm)      media 1993-2019 (mm)
## Length:13                      Min.   : 7.0      Min.   : 57.0
## Class :character               1st Qu.: 28.0    1st Qu.: 74.0
## Mode  :character               Median : 89.0    Median : 96.0
##                                Mean    : 180.2    Mean    : 171.5
##                                3rd Qu.: 171.0    3rd Qu.: 114.0
##                                Max.     : 1171.0   Max.     : 1114.0
```

Esercizio: Importate i dati inserendo i parametri corretti di read.csv - come noterete, la prima riga NON importa correttamente. La seconda sì, ma è una funzione della libreria data.table. Usate la documentazione in R per capire come aggiungere argomenti alla funzione read.csv (vedi sezione QUI). **NB** usate il tasto TAB per richiamare i percorsi, risparmiate tempo e fatica!

## 2.3 Input dati: Fogli excel

I fogli di calcolo in formato MS Excel sono più strutturati rispetto ad un file CSV, dato che ogni colonna ha anche lo specifico tipo di dati (numero, testo, data ecc...).

Per leggere/scrivere questi dati è necessario usare librerie aggiuntive.

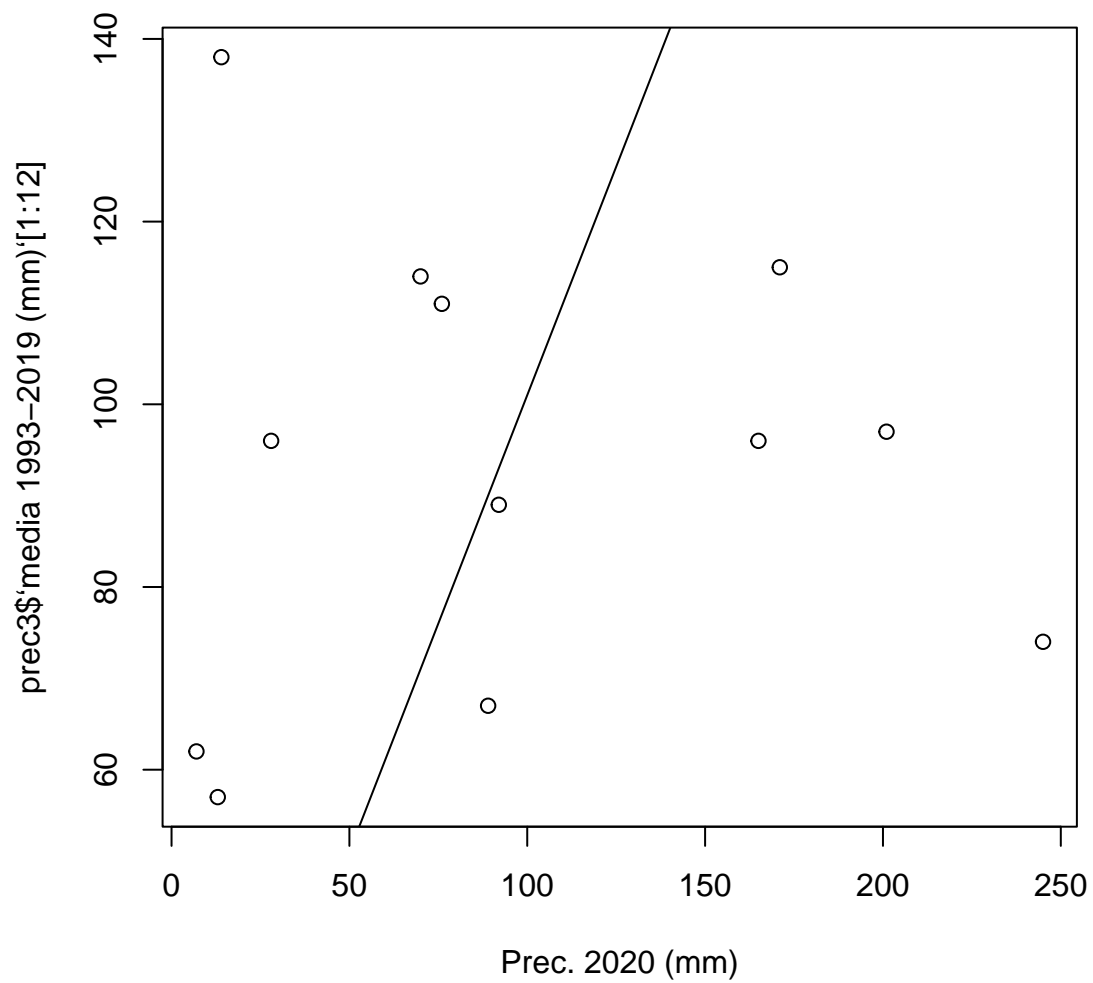
```
prec3 <- prec2
## nuova colonna
prec3$norm2020 <- prec3$`2020 (mm)` / prec3$`media 1993-2019 (mm)`
writexl::write_xlsx( list("originale" =prec2, "elaborata"=prec3), path="dati/output.xlsx" )
prec4 <- readxl::read_xlsx("dati/output.xlsx")
```

	A	B	C	D	E
1	Precipitazioni mensili Veneto	2020 (mm)	media 1993-2019 (mm)	norm2020	
2	Gennaio	13	57	0.228070175	
3	Febbraio	7	62	0.112903226	
4	Marzo	89	67	1.328358209	
5	Aprile	28	96	0.291666667	
6	Maggio	70	114	0.614035088	
7	Giugno	165	96	1.71875	
8	Luglio	92	89	1.033707865	
9	Agosto	201	97	2.072164948	
10	Settembre	76	111	0.684684685	
11	Ottobre	171	115	1.486956522	
12	Novembre	14	138	0.101449275	
13	Dicembre	245	74	3.310810811	

## 2.4 Grafici di base

La funzione plot riesce a creare grafici di base. Andiamo a confrontare le precipitazioni del 2020 con quelle medie 1993-2019.

```
plot(prec3$`2020 (mm)`[1:12], prec3$`media 1993-2019 (mm)`[1:12], xlab="Prec. 2020 (mm)")
abline(1,1)
```

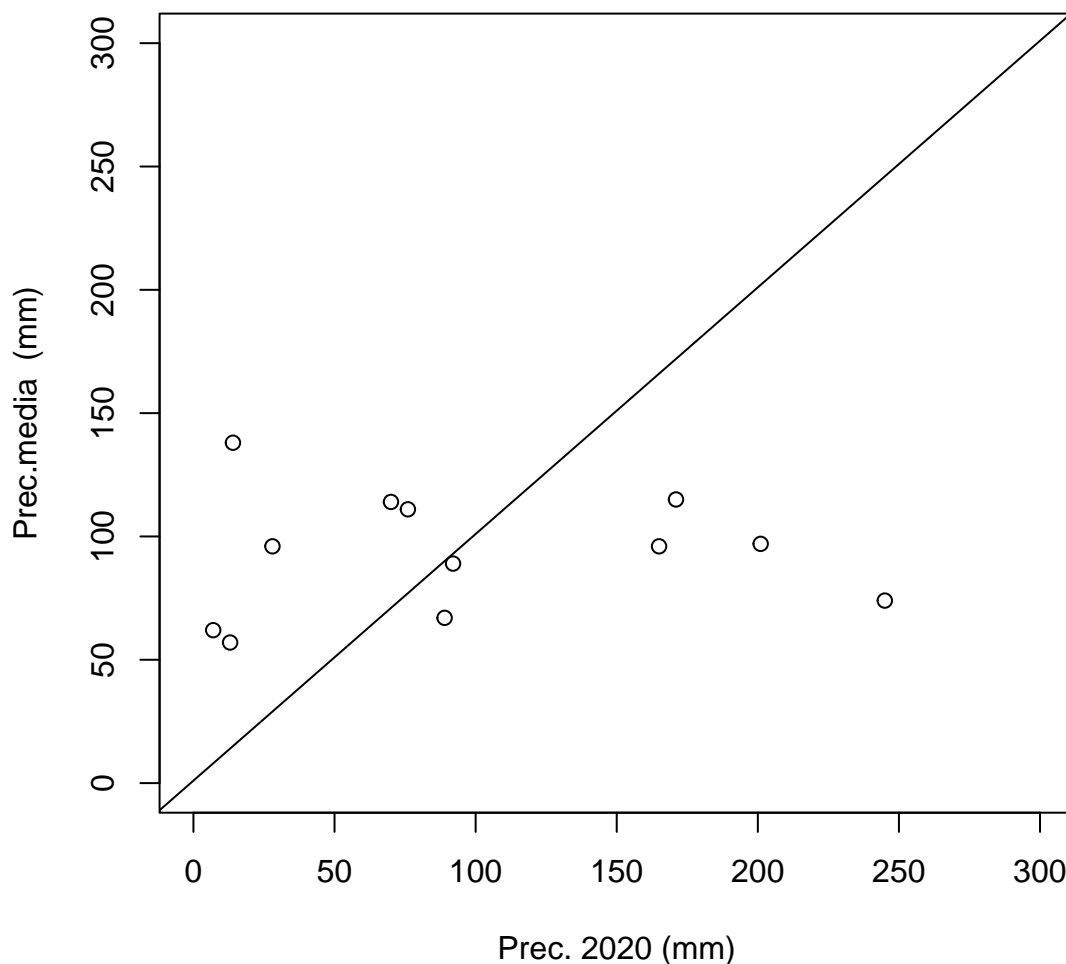


## 2.5 Grafici di base

Miglioriamo un po' il grafico - ma utilizzeremo poi la libreria ggplot2 per ottenere grafici migliori.

```
plot(prec3$`2020 (mm)`[1:12], prec3$`media 1993-2019 (mm)`[1:12],
     xlab="Prec. 2020 (mm)", ylab="Prec.media (mm)",
     xlim=c(0,300), ylim=c(0,300))
abline(1,1)
```





## 2.6 Input dati: JSON

I dati JSON (**J**ava**S**cript **O**bject **N**otation) sono molto utilizzati per trasferire dati non strutturati. Vedi QUI e QUI. Anche per dati “geo” con il GeoJSON (vedi esempio QUI).

Usiamo direttamente da API ARPAV “Anagrafica stazioni in formato JSON” dalla pagina QUI

```
## Scarichiamo direttamente da API ARPAV
api.url <- "https://api.arpa.veneto.it/REST/v1/meteo_meteogrammi?rete=MGRAMMI&coordcd=18&orario=0&rn
stazioni.meteo <- jsonlite::read_json(api.url)
```

## 3 Trasformazione dei dati

### 3.1 Introduzione

Modificare dati significa aggiungere, togliere, filtrare, manipolare colonne, creare sommari e statistiche di interesse aggregando dati. Insomma il dato viene “manipolato”. In inglese viene anche usato “data wrangling”

Ci “alleniamo” con i dati JSON delle stazioni meteo ARPAV scaricati dal JSON.

**Obiettivo:** creare un grafico con la posizione delle stazioni e aggiungere qualche valore mediante unione con altri dati.

## 3.2 Iteratori

Fondamentale saper utilizzare operazioni che eseguono su tutti gli elementi di un oggetto. E.g. tutte le colonne di un *data.frame*, tutte le righe di un *data.frame*, tutti gli oggetti di una *list*...

```
dati <- stazioni.meteo$data
length(dati)
## LOOP classico: visto male per
## performance bassa,
## ma non necessariamente vero
for(i in 1:length(dati)){
  print(dati[[i]])
  break # esce dal ciclo
  next # continua il ciclo senza elaborare le righe successive
  print("Questa riga non viene elaborata")
}

funz <- function(x){
  length(x)
}
# lapply/mapply/sapply/apply per righe/colonne
# togliete il commento alla riga sotto
# lapply(dati, funz)
```

## 4 Tidyverse

### 4.1 Introduzione

Tidyverse, dal sito: “un set di funzioni e librerie di R per la *data science*. Tutti i pacchetti condividono una filosofia di progettazione, una grammatica e strutture di dati”.

La grammatica per la manipolazione dei dati fornisce un insieme coerente di funzioni (verbi) dalla libreria *dplyr*:

- **mutate()** aggiunge nuove variabili che sono funzioni di variabili esistenti
- **select()** seleziona le variabili in base al loro nome.
- **filter()** seleziona i casi in base ai loro valori.
- **summarise()** riduce più valori a un singolo riassunto.
- **arrange()** cambia l'ordine delle righe.

Combinando con **group\_by()** questa operazione viene stratificata “per gruppo” aggregando le operazioni. Ci sono molti altri operatori ma questi sono quelli di base.

### 4.2 Iteratori: map

L'operatore “map” è simile agli “apply” ma con più funzionalità. Usando *map\_df* i risultati vengono raccolti su una tabella (tibble)!

```
library(tidyverse)
## i dati sono strutturati nell'oggetto "data"
## come liste dentro altre liste

dt <- map(.x = stazioni.meteo$data, .f = function(x){ x } )
dtf <- map_df (.x = stazioni.meteo$data, .f = function(x){ x } )
```

### 4.3 Iteratori: map

Mappando a data.frame con `map` abbiamo colonne di tipo character che dovrebbero essere numeric dobbiamo in alcuni casi assegnare il tipo alla colonna come sotto.

NB gli iteratori su tabelle (data.frame o tibble) lavorano per colonne (column-wise).

```
## richiamo la classe di ogni colonna
# map(.x = dt, .f = function(x){ class(x) } )
map(.x = dtf, .f = function(x){ class(x) } )
## ATTENZIONE - colonna valore è corretta?
dtf$valore <- as.numeric(dtf$valore)
dtf$tipo <- as.factor(dtf$tipo)
levels(dtf$tipo)
```

### 4.4 Iteratori: map

NOTA che `map` raccoglie ogni elemento della lista in un vettore ed il vettore può avere solo valori di un tipo base (*character, numeric, integer, logical...*) - sceglie sempre il *character* se trova diverse tipologie.

NB con `map_df` questo non succede perchè la funzione `map_df` è “smart” e capisce qualche tipo usare per ogni colonna. TRANNE che per la colonna valore. Perché?\*

\*\* Nota che nel JSON il numero in “valore” è tra virgolette, dunque viene visto come testo (character) e non come numerico. Confronta con latitudine e longitudine!

```
str(stazioni.meteo$data[[1]])
```

```
## List of 15
## $ codice_stazione: int 19
## $ codseqst       : int 300000170
## $ nome_stazione  : chr "Agordo"
## $ point          : chr "{\"type\": \"Point\", \"coordinates\": [12.03319083, 46.27804651]}\"
## $ longitudine    : num 12
## $ latitudine     : num 46.3
## $ quota          : num 585
## $ nome_sensore   : chr "Temperatura aria a 2m"
## $ dataora        : chr "2024-09-05T19:00:00"
## $ valore         : chr "16.0"
## $ misura         : chr "°C"
## $ tipo           : chr "TARIA2M"
## $ aggiornamento  : chr "2024-09-05T19:15:06"
## $ gestore        : chr "ARPAV"
## $ provincia      : chr "BELLUNO"
```

### 4.5 Filter e select

Operazioni usando tidyverse (e altri sistemi) usano operatori di concatenazione di processi. In questo caso si usa `%>%` ovvero un operatore che indica di usare l’output di un processo come input al processo successivo. Un esempio con un’operazione di filtro con successiva selezione di colonne:

```
## Selezioniamo solo le stazioni in provincia di Belluno
dtf.belluno <- dtf %>% filter( provincia == "BELLUNO")
## Se aggiungiamo criteri di filtro divisi da virgola è come usare AND
dtf %>% filter( provincia != "BELLUNO", valore < 10)
dtf %>% filter( provincia != "BELLUNO" & valore < 10)
## Se aggiungiamo criteri di filtro è come usare AND
dtf %>% filter( provincia == "BELLUNO" | valore < 12) %>% select(nome_stazione, valore)
```

## 4.6 Mutate

Il comando mutate consente di calcolare nuove colonne in funzione di quelle esistenti.

```
dtf %>% mutate(valore=valore^2) %>% select(valore)
```

## 4.7 Aggregazioni

Operazioni di aggregazione con `group_by` e `summarise`

Attenzione - la funzione `mean` restituisce NA se tra i valori trova degli NA, per evitare questo comportamento mettere `na.rm = TRUE`.

```
dtf %>% group_by(tipo) %>% summarise(valore=mean(valore))
```

```
## # A tibble: 3 x 2
##   tipo    valore
##   <fct>    <dbl>
## 1 TARIA2M    NA
## 2 TARIA5M    NA
## 3 TARIA8M   9.66
```

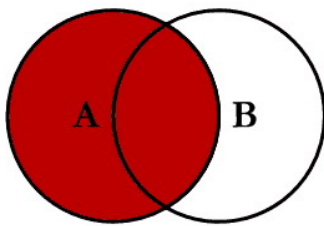
```
dtf %>% group_by(tipo) %>% summarise(valore=mean(valore, na.rm = TRUE))
```

```
## # A tibble: 3 x 2
##   tipo    valore
##   <fct>    <dbl>
## 1 TARIA2M  17.1
## 2 TARIA5M   8.88
## 3 TARIA8M   9.66
```

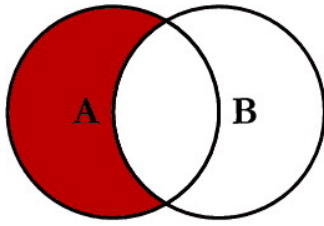
## 4.8 Unione tabelle (join)

I join sono operazioni note in ambito di gestione dei dati, servono per aggiungere in senso orizzontale valori usando una o più colonne con valori corrispondenti.

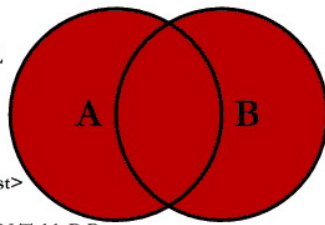
# SQL JOINS



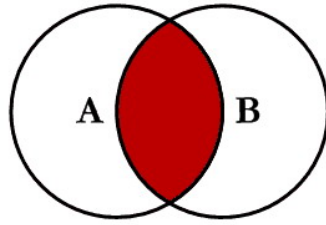
```
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
```



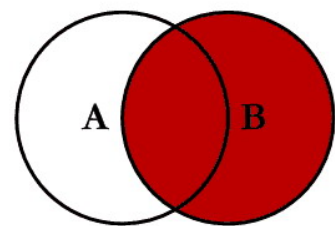
```
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
WHERE B.Key IS NULL
```



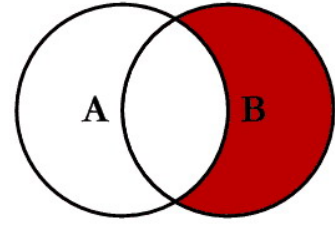
```
SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key
```



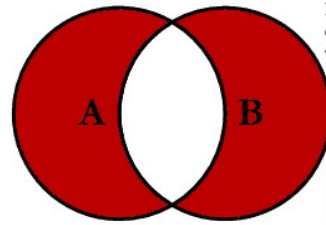
```
SELECT <select_list>
FROM TableA A
INNER JOIN TableB B
ON A.Key = B.Key
```



```
SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
```



```
SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL
```



```
SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL
OR B.Key IS NULL
```

© C.L. Moffatt, 2008

## SQL Joins Cheatsheet

[ml4devs.com/sql-joins](http://ml4devs.com/sql-joins)



INNER JOIN

```
SELECT <columns>
FROM A
[INNER] JOIN B
ON A.Key = B.Key;
```



LEFT JOIN

```
SELECT <columns>
FROM A
LEFT [OUTER] JOIN B
ON A.Key = B.Key;
```



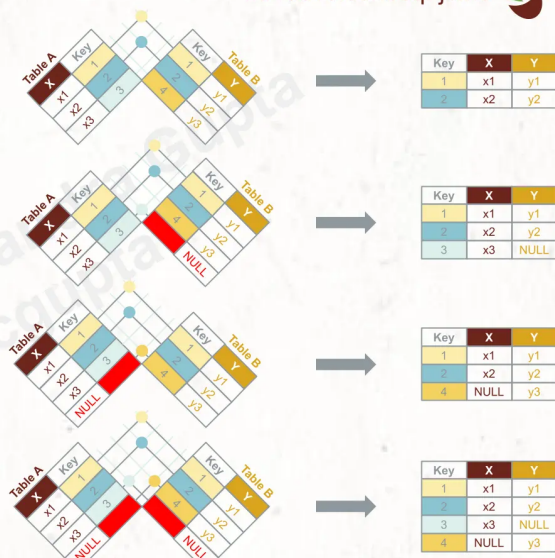
RIGHT JOIN

```
SELECT <columns>
FROM A
RIGHT [OUTER] JOIN B
ON A.Key = B.Key;
```



FULL JOIN

```
SELECT <columns>
FROM A
FULL [OUTER] JOIN B
ON A.Key = B.Key;
```



© Satish Chandra Gupta, Creative Commons BY-NC-ND 4.0 International License.

scgupta

[linkedin.com/in/scgupta](https://www.linkedin.com/in/scgupta)

## 4.9 Unione tabelle (join)

Vediamo un esempio pratico: andiamo ad unire i valori di oggi a quelli di alcuni giorni fa alle nostre stazioni. Scarica qui il dato CSV.

```
dt.past <- read.csv("dati/ARPA_temperatura_20240528.csv")
dt.past <- dt.past %>% select(statnm, statcd, dataora, media)
dt.joined <- dt.past %>% left_join(dtf, by=c("statcd"="codice_stazione") )
```

## 4.10 Esercizio 2

- Esportare un foglio excel ed un file CSV separato da TAB con i dati di precipitazione in Veneto fatti fino a questo punto (df.joined) ed un file di R che li legge in 2 oggetti tibble.
- consegnate: i seguenti file:

- es2\_cognome\_nome.xlsx
- es2\_cognome\_nome.csv
- es2\_cognome\_nome.R

## 4.11 Ristrutturazione: wide=>long

Questo tipo di modifica è molto comune per rappresentare i dati in modo differente. Vediamo un esempio pratico:

```
## andiamo a creare classi di temperatura
classi <- cut(dtf$valore, breaks = c(-10,0,10,20,30))
## mappatura con tabella a doppia entrata
tb.wide <- table(dtf$provincia, classi)
## nota che è un oggetto "table", trasformiamo in data.frame
tb.wide <- as.data.frame.matrix(tb.wide)
rownames(tb.wide)
```

```
## [1] "BELLUNO" "PADOVA" "ROVIGO" "TREVISO" "UDINE" "VENEZIA" "VERONA"
## [8] "VICENZA"
```

```
## nota ancora che le provincie sono rownames - se converto
## a tibble non ci sono più
rownames(as_tibble(tb.wide) )
```

```
## [1] "1" "2" "3" "4" "5" "6" "7" "8"
```

```
## creo colonna e relocate la mette all'inizio (provate senza "relocate")
tb.wide <- tb.wide %>% mutate(Provincie = rownames(tb.wide))%>% relocate(Provincie)
## tabella con numero di stazioni per classe espressa "wide"
tb.wide
```

```
##      Provincie (-10,0] (0,10] (10,20] (20,30]
## BELLUNO    BELLUNO      0      17      34      0
## PADOVA      PADOVA      0       0       8       2
## ROVIGO      ROVIGO      0       0       4       2
## TREVISO     TREVISO      0       0      13       2
## UDINE       UDINE      0       0       1       0
## VENEZIA     VENEZIA      0       0      10       5
## VERONA      VERONA      0       0       8       5
## VICENZA     VICENZA      0       1      32       2
```

```
## trasformo
tb.long <- tb.wide %>% pivot_longer(!Provincie,
  names_to = "Classe Temp. °C", values_to = "Numero"
)
tb.long
```

```
## # A tibble: 32 x 3
##   Provincie `Classe Temp. °C` Numero
##   <chr>      <chr>          <int>
## 1 BELLUNO  (-10,0]                0
## 2 BELLUNO   (0,10]               17
## 3 BELLUNO  (10,20]              34
## 4 BELLUNO  (20,30]                0
## 5 PADOVA   (-10,0]                0
## 6 PADOVA    (0,10]                0
## 7 PADOVA   (10,20]                8
## 8 PADOVA   (20,30]                2
## 9 ROVIGO   (-10,0]                0
## 10 ROVIGO  (0,10]                0
## # i 22 more rows
```

## 4.12 Ristrutturazione: long=>wide

Operazione inversa per tornare alla tabella di origine

```
tb.long2wide <- tb.long %>% pivot_wider(
  names_from = "Classe Temp. °C", values_from = "Numero"
)
tb.long2wide
```

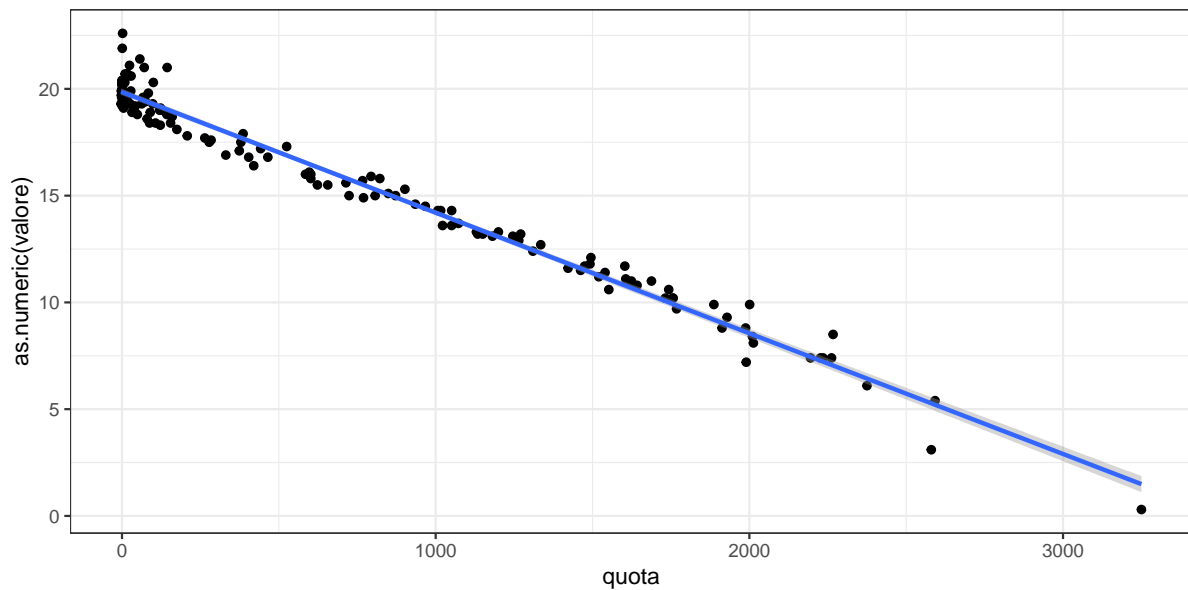
```
## # A tibble: 8 x 5
##   Provincie `(-10,0]` `(0,10]` `(10,20]` `(20,30]`
##   <chr>      <int>    <int>    <int>    <int>
## 1 BELLUNO      0      17      34      0
## 2 PADOVA       0       0       8       2
## 3 ROVIGO       0       0       4       2
## 4 TREVISO      0       0      13       2
## 5 UDINE        0       0       1       0
## 6 VENEZIA      0       0      10       5
## 7 VERONA       0       0       8       5
## 8 VICENZA      0       1      32       2
```

## 5 Grafici

### 5.1 Grafici con GGplot - dispersione

Le funzioni di base “plot” sono utili ma ci sono librerie che forniscono più funzionalità. Introduciamo “ggplot”. Notate che viene creato un grafico concatenando con l’operatore “+” gli oggetti del grafico. “aes” sta per *aesthetics* e indica colori e valori degli assi, aggregazioni ecc...

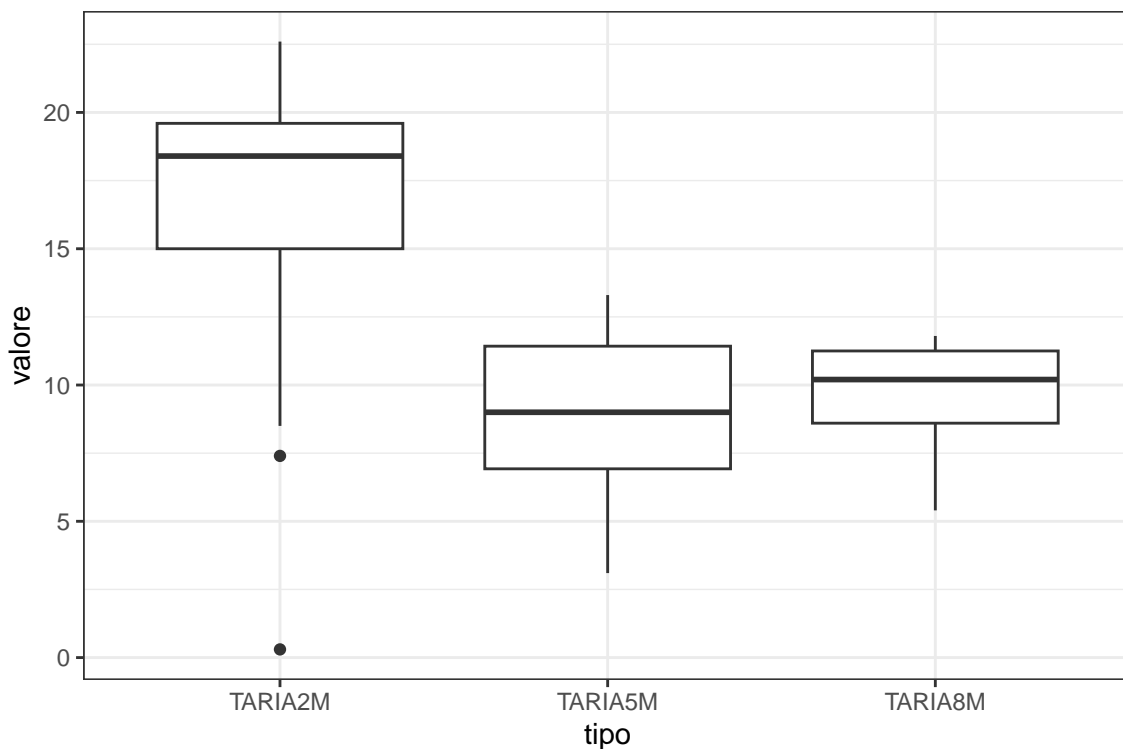
```
library(ggplot2)
ggplot(dtf, aes(x=quota, y=as.numeric(valore))) + geom_point( ) +
  geom_smooth(method = "lm")+ theme_bw()
```



## 5.2 Grafici con GGplot - scatola e baffi

NB. la parte “as.numeric” in quanto la colonna anche contiene numeri è di tipo “character”. Modifichiamola con “as.numeric”, ma vedrete due valori mancanti (NA=not available).

```
dtf$valore <- as.numeric(dtf$valore)
ggplot(dtf) + geom_boxplot(aes(x=tipo, y=valore)) + theme_bw()
```

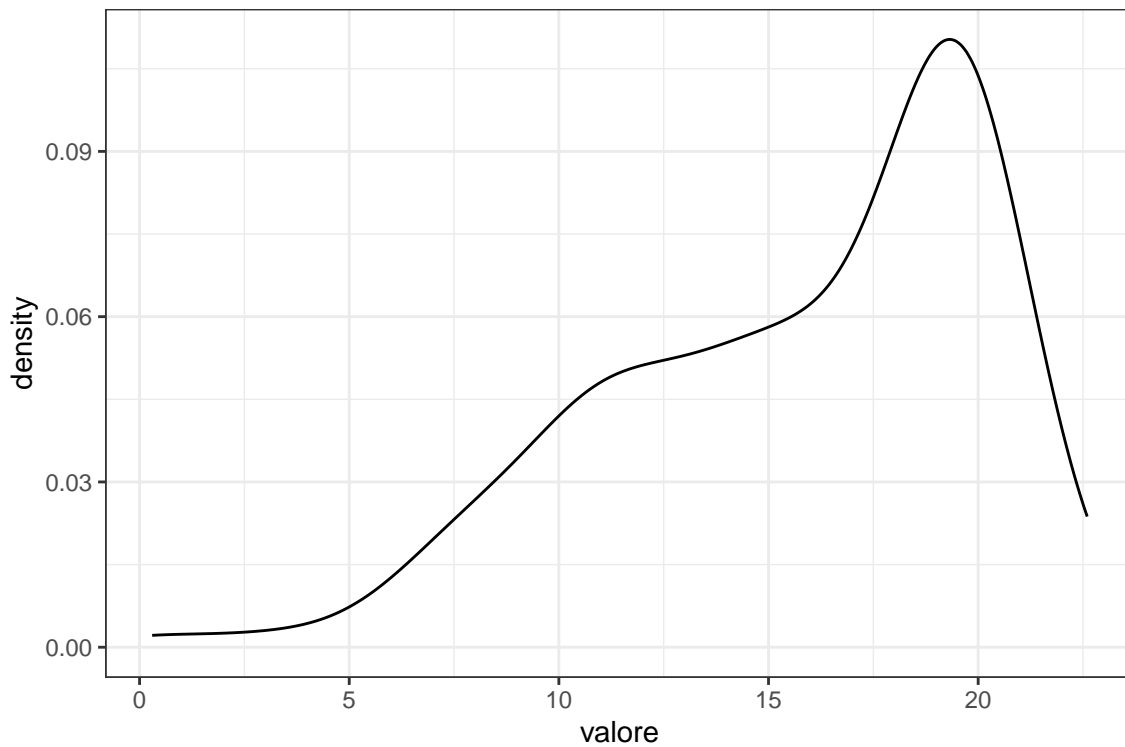


## 5.3 Grafici multipli (facets)

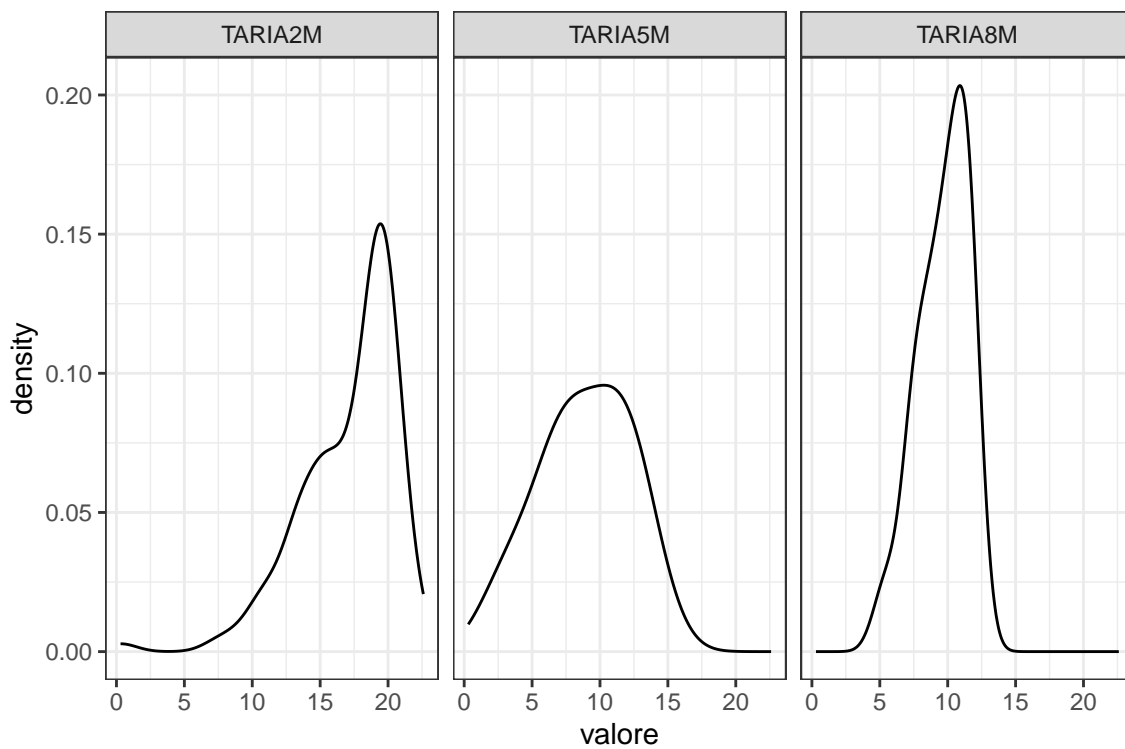
`facet_wrap` e `facet_grid` consentono di “esplodere” per una o più variabili di raggruppamento. Qui un esempio con una stima della funzione di distribuzione delle frequenze.



```
ggplot(dtf) + geom_density(aes(x=valore )) + theme_bw()
```



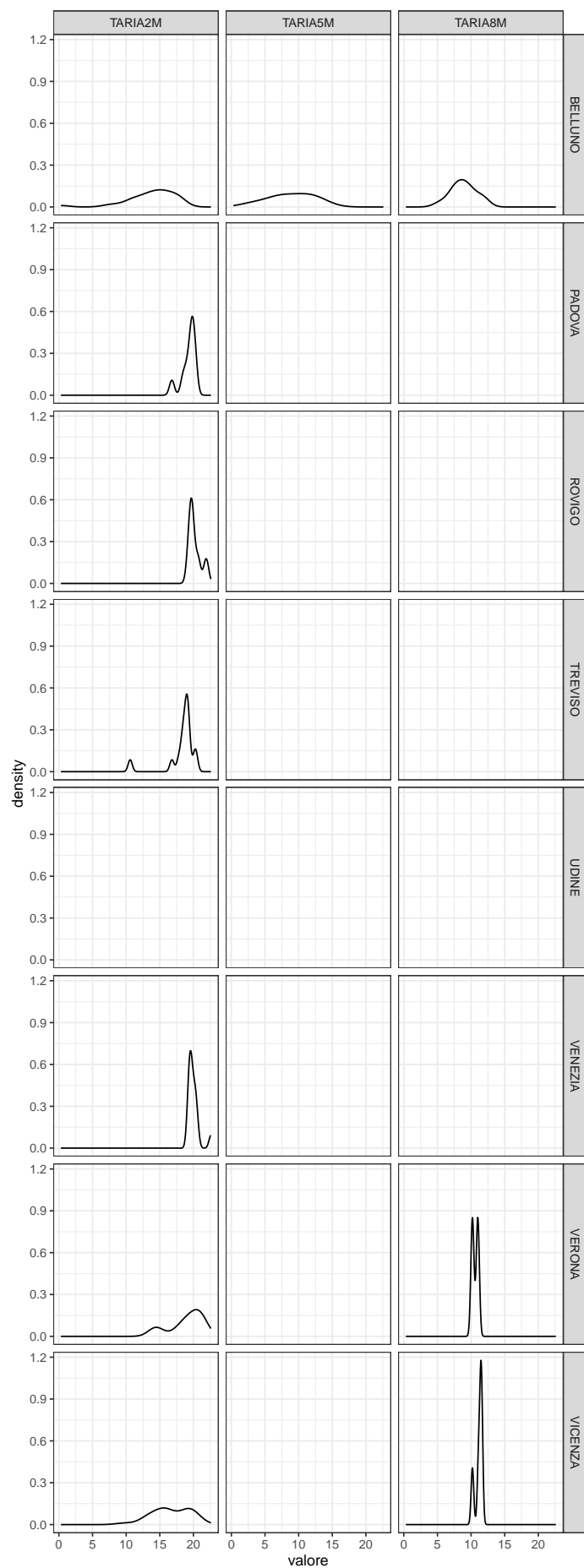
```
ggplot(dtf) + geom_density(aes(x=valore)) + facet_wrap(vars(tipo)) + theme_bw()
```



## 5.4 Grafici multipli (facets)

`facet_wrap` e `facet_grid` consentono di “esplodere” per una o più variabili di raggruppamento.

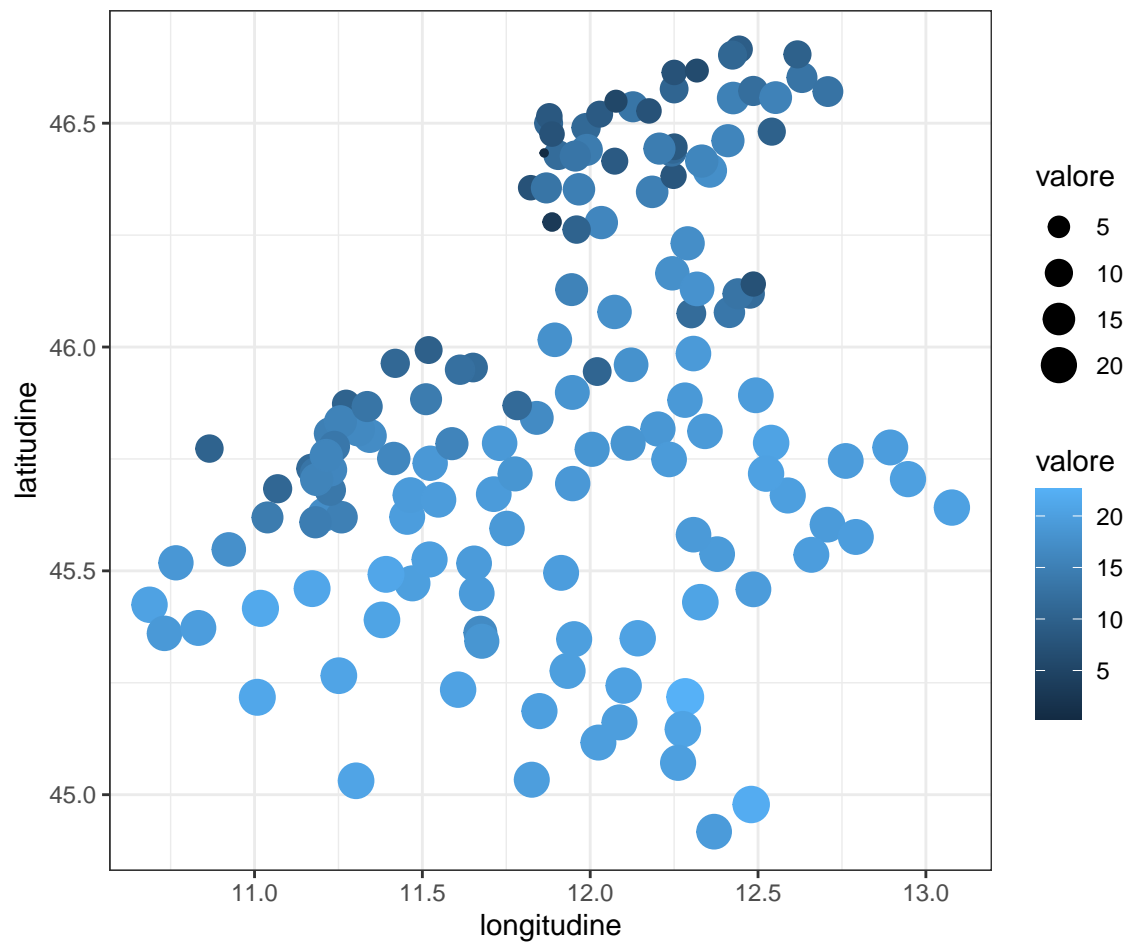
```
ggplot(dtf) + geom_density(aes(x=valore)) + facet_grid(provincia~tipo ) + theme_bw()
```



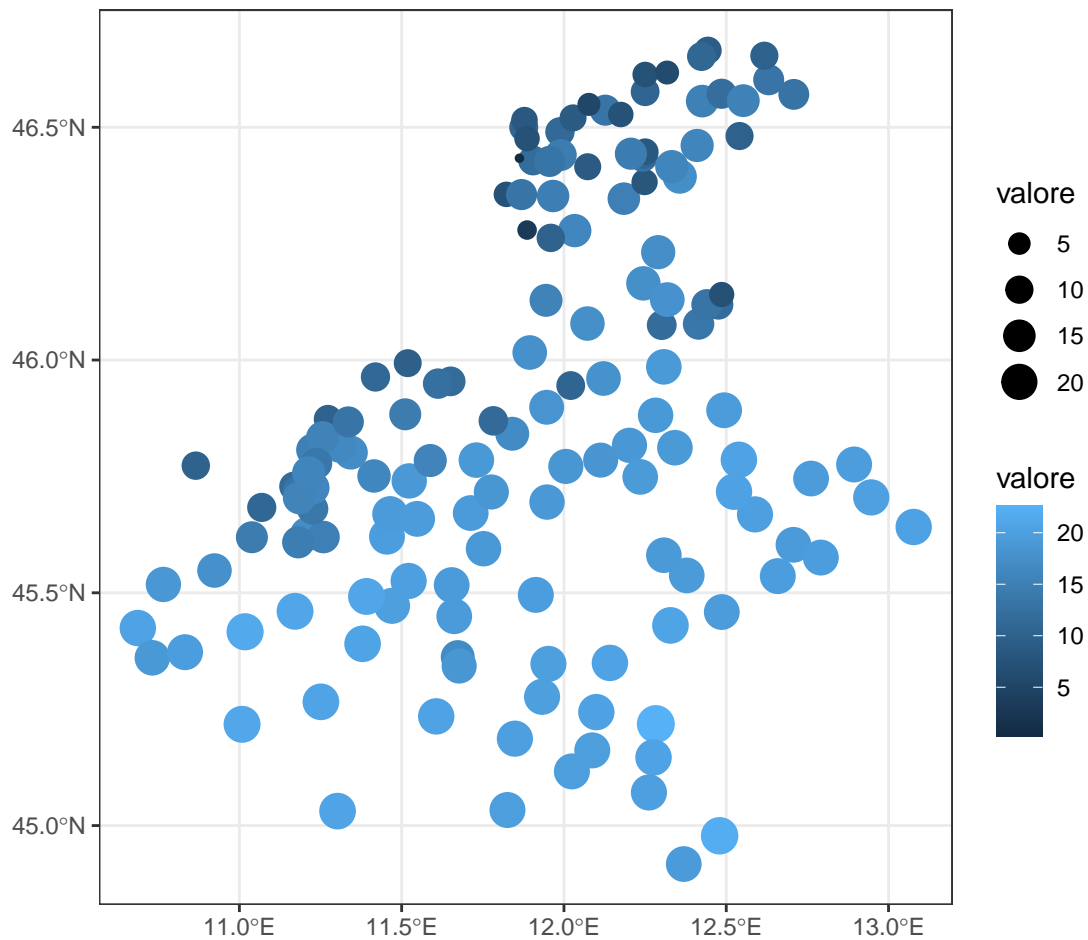
## 5.5 Grafici e dati geo

Abbiamo colonne con coordinate: possiamo creare un grafico con uno sfondo territoriale e considerare delle colonne come coordinate? Per dati geo usiamo libreria *sf* e *terra*. Non indispensabili per geostatistica, ma fortemente consigliati.

```
library(ggspatial, quietly = T)
library(sf)  ## se non l'avete installata non vi preoccupate
theme_set(theme_bw())
ggplot(dtf) + geom_point(aes(x=longitudine, y=latitudine, color=valore, size=valore ))
```



```
dtf.geo <- sf::st_as_sf(dtf, coords = c("longitudine", "latitudine"), crs=4326 )
ggplot(dtf.geo) + geom_sf(aes(color=valore, size=valore )) + coord_sf()
```



## 5.6 Grafici e stampe PNG

Il metodo per avere grafici su stampa mediante immagini o pdf

```
png("fileGrafico.png")

p <- ggplot(dtf.geo) + geom_sf(aes(color=valore, size=valore )) + coord_sf()
print(p)
dev.off()
```

possiamo impostare dimensioni in pixel (width e height) ed inoltre la risoluzione (res)... provate a vedere.

```
png("fileGrafico300dpi.png", width=2000, height=2000, res=300)
p <- ggplot(dtf.geo) + geom_sf(aes(color=valore, size=valore )) + coord_sf()
print(p)
dev.off()
```

## 5.7 Grafici e stampe PDF

Il formato PDF ha il vantaggio/svantaggio di essere vettoriale (non perde definizione) - provate a vedere la definizione. Perché non usiamo SEMPRE pdf? Immaginatevi a inserire milioni di punti... perde efficacia.

```
pdf("fileGrafico.pdf")
p <- ggplot(dtf.geo) + geom_sf(aes(color=valore, size=valore )) + coord_sf()
print(p)
dev.off()
```

NB possiamo creare un iterazione per fare un grafico per foglio in un singolo file PDF!

```
plotMe <- function(tipoStazione){
  dati <- dtf.geo %>% filter(tipo==tipoStazione)
  p <- ggplot(dati) + geom_sf(aes(color=valore)) + ggtitle(paste("Tipo: ", tipoStazione)) + coord_sf()
  print(p)
}
## NB la colonna "tipo" è di tipo FACTOR
pdf("fileGraficoMulti.pdf")
for(tipoStz in levels(dtf.geo$tipo) ){
  p <- plotMe(tipoStz)
  print(p)
}
dev.off()
```

## 5.8 Cheatsheets

Sono PDF formato poster con indicazioni di base sull'utilizzo di strumenti R

# Data transformation with dplyr : : CHEATSHEET



**dplyr** functions work with pipes and expect **tidy data**. In tidy data:

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- x** > **f(y)** becomes **f(x, y)**

## Summarize Cases

Apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

**summary function**

- summarize(data, ...)** Compute table of summaries.  
mtcars > summarize(avg = mean(mpg))
- count(data, ..., wt = NULL, sort = FALSE, name = NULL)** Count number of rows in each group defined by the variables in ... Also **tally()**, **add\_count()**, **add\_tally()**.  
mtcars > count(cyl)

## Group Cases

Use **group\_by(data, ..., add = FALSE, drop = TRUE)** to create a "grouped" copy of a table grouped by columns in ... dplyr

## Manipulate Cases

### EXTRACT CASES

Row functions return a subset of rows as a new table.

- filter(data, ...)** Extract rows that meet logical criteria.  
mtcars > filter(mpg > 20)
- distinct(data, ..., keep\_all = FALSE)** Remove rows with duplicate values.  
mtcars > distinct(gear)
- slice(data, ..., preserve = FALSE)** Select rows by position.  
mtcars > slice(10:15)
- slice\_sample(data, ..., n, prop, weight\_by = NULL, replace = FALSE)** Randomly select rows. Use **n** to select a number of rows and **prop** to select a fraction of rows.  
mtcars > slice\_sample(n = 5, replace = TRUE)
- slice\_min(data, order\_by, ..., n, prop, with\_ties = TRUE)** and **slice\_max()** Select rows with the lowest and highest values.  
mtcars > slice\_min(mpg, prop = 0.25)
- slice\_head(data, ..., n, prop)** and **slice\_tail()** Select the first or last rows.  
mtcars > slice\_head(n = 5)

## Manipulate Variables

### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

- pull(data, var = 1, name = NULL, ...)** Extract column values as a vector, by name or index.  
mtcars > pull(wt)
- select(data, ...)** Extract columns as a table.  
mtcars > select(mpg, wt)
- relocate(data, ..., before = NULL, after = NULL)** Move columns to new position.  
mtcars > relocate(mpg, cyl, after = last\_col())

**Use these helpers with select() and across()**  
e.g. mtcars > select(mpg:cyl)

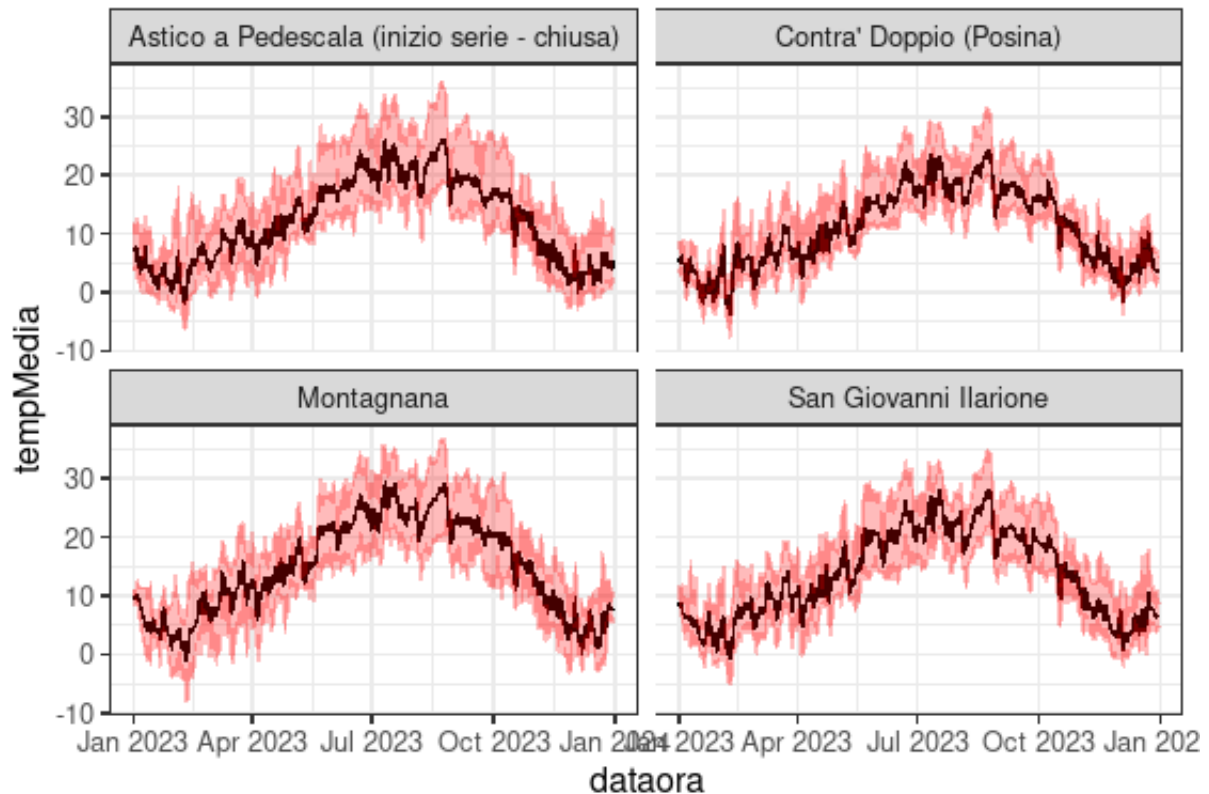
<b>contains(match)</b>	<b>num_range(prefix, range)</b>	! e.g., mpg:cyl
<b>ends_with(match)</b>	<b>all_of(x)/any_of(x, ..., vars)</b>	! e.g., !gear
<b>starts_with(match)</b>	<b>matches(match)</b>	<b>everything()</b>

### MANIPULATE MULTIPLE VARIABLES AT ONCE

**across(cols, funs, ..., names = NULL)** Summarize or mutate multiple columns in the same way.

## 5.9 Esercizio 3

Vedi file "modulo1\_03a\_esercizio\_finale2.R" e "modulo1\_03b\_esercizio\_finale2.R" nella cartella "esercizi/modulo1" - eseguite il seguente grafico estraendo i dati dalle API meteo ARPAV.



## 6 Dati geospaziali

### 6.1 Introduzione

A grandi linee i dati spaziali sono interpretati digitalmente mediante due modelli: il modello vettoriale ed il modello raster. In entrambi i casi va comunque compreso che ogni elemento è georiferito, ovvero ha una localizzazione in un sistema di riferimento, quindi delle coordinate.

### 6.2 Dati vettoriali

I dati vettoriali hanno tre primitive geometriche di base: punti linee e poligoni. La base è sempre il punto con una coordinate, in 2d o 3d.

Ci sono decine di formati di file che rappresentano dati vettoriali. Ne abbiamo visto uno all'inizio, il file JSON! Per un riferimento vedete la pagina [QUI](#). Ad oggi la libreria di riferimento per leggere questi dati è la libreria `sf`.

### 6.3 Raster

Il formato raster è molto semplice: sono nodi a passo costante, quindi una “griglia” georiferita mediante uno o più punti. Il passo della griglia è spesso identico nelle due direzioni (ma non necessariamente), e ne determina la risoluzione spaziale.

Ci sono decine di formati di file che rappresentano dati raster. Per un riferimento vedete la pagina [QUI](#). Ad oggi la libreria di riferimento per leggere questi dati è la libreria `terra`.

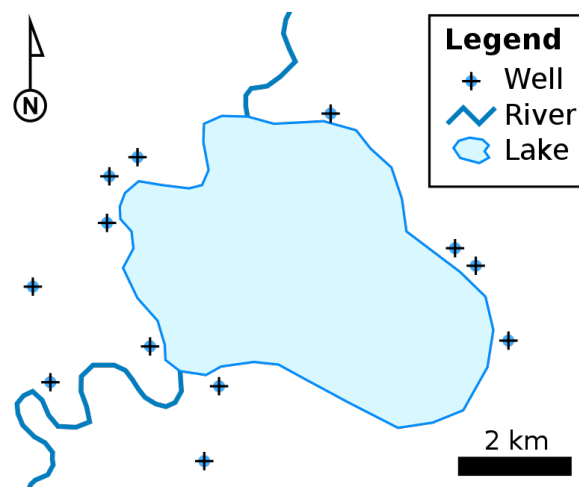


Figure 1: By M. W. Toews - Own work, CC BY-SA 4.0

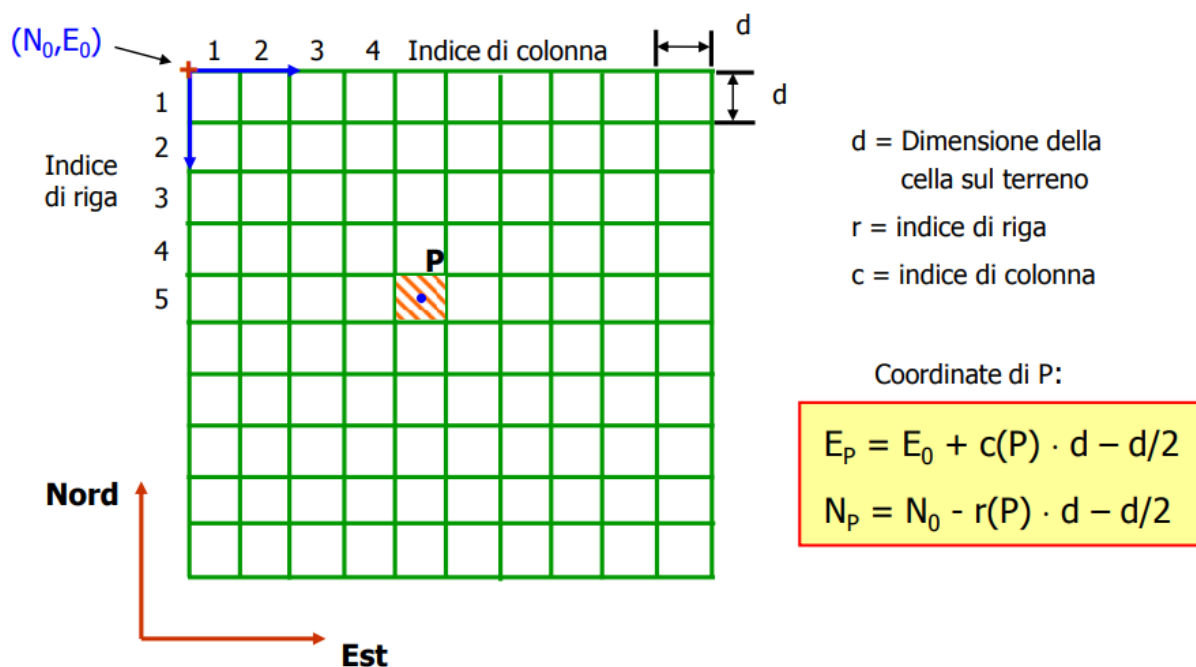


Figure 2: fonte: Appunti GIS Prof. Alberto Guarnieri



## 6.4 Interpolazione

Stimare valori di una variabile applicando un modello creato conoscendo valori misurati (e.g. campionamenti). E' ragionevole adottare un criterio spaziale, ovvero usare un criterio di spazio (il valore del punto sarà probabilmente più simile ad una misura vicina che una misura più lontana, a parità di altre condizioni). Da qui la geostatistica.

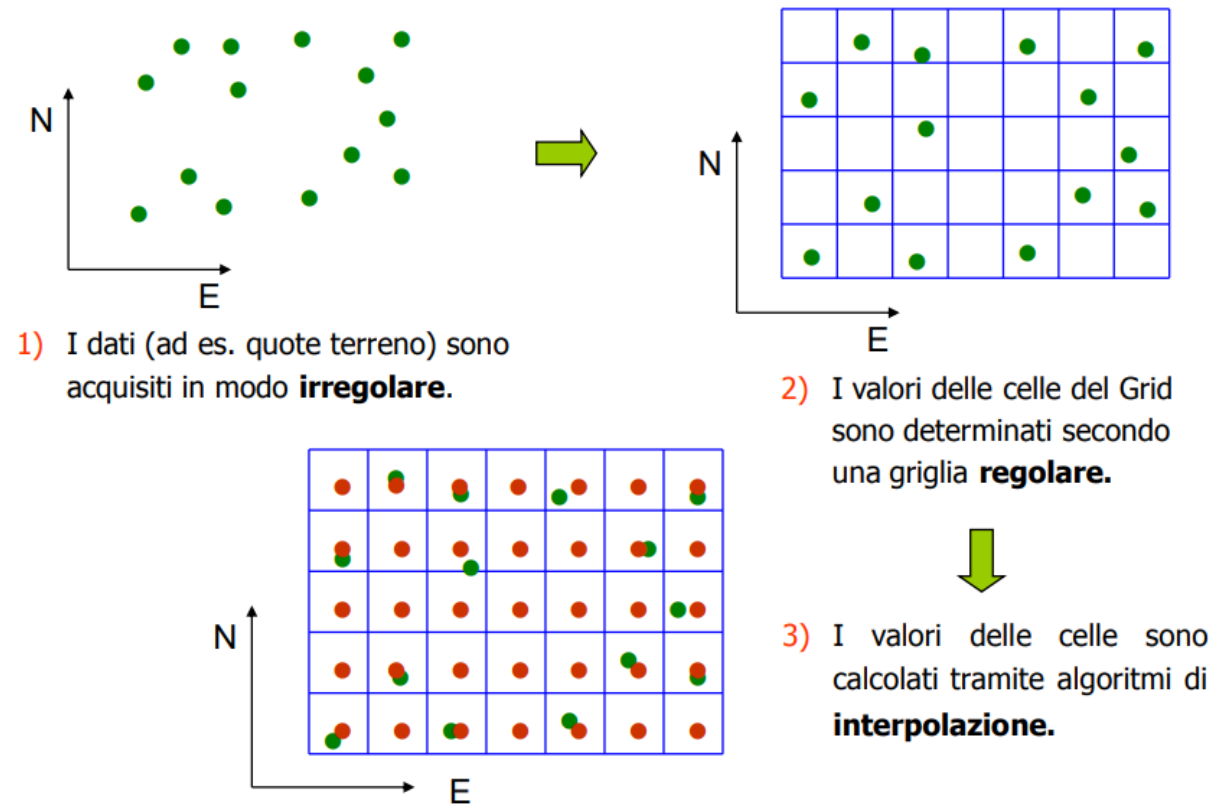


Figure 3: fonte: Appunti GIS Prof. Alberto Guarnieri