

UNIVERSITÀ DEGLI STUDI DI MILANO-BICOCCA

Scuola di Economia e Statistica

Corso di laurea in Scienze Statistiche ed Economiche



**IL CORONAVIRUS: MODELLI STATISTICI
PER L'ANALISI DELLA DIFFUSIONE
NEL TERRITORIO LOMBARDO**

Relatore:

Ch.mo prof. Aldo Solari

Tesi di Laurea di:

Matteo Baldanza

Matr. N. 826018

ANNO ACCADEMICO

2019/2020

Ai miei genitori e
alla mia famiglia

*“Se torturi i numeri a lungo,
confesseranno qualsiasi cosa”
Gregg Easterbrook*

Indice

Elenco delle figure	5
Introduzione	6
1 Analisi dei dati	8
1.1 L'analisi descrittiva dei dati	8
1.2 La prevalenza: Italia e Hubei	10
1.3 Incidenza: Italia e Hubei	12
1.4 Analisi delle varie misure e problematiche annesse	14
1.5 La mortalità in Italia	15
2 La stima di R_0	19
2.1 Introduzione ai modelli epidemiologici	19
2.2 Il parametro R_0	20
2.2.1 Alcune delucidazioni su R_0 e differenze con R_t	22
2.3 Generation-Time	23
2.4 Stima di R_0 tramite il tasso di crescita esponenziale	25
2.4.1 Introduzione	25
2.4.2 L'equazione di Lotka-Eulero	26
2.4.3 Derivazione di R_0 dall'equazione di Lotka-Eulero	28
2.4.4 R_0 utilizzando campione di osservazioni cinesi e italiane	29
2.5 Stima di R_0 tramite la funzione di massima verosimiglianza	31
2.5.1 Derivazione della stima di massima verosimiglianza	31
2.5.2 R_0 utilizzando campione di osservazioni cinesi e italiane	35

3	I modelli epidemiologici	36
3.1	Modello <i>SIR</i>	36
3.1.1	In termini teorici	36
3.1.2	Relazione R_0 con il numero di infetti	38
3.1.3	Applicazione modello SIR in Lombardia	39
3.1.4	L'andamento di R_t in Lombardia	41
3.2	Modello <i>SEIR</i>	45
3.2.1	In termini teorici	45
3.2.2	Applicazione del modello in Lombardia	47
4	La struttura dei mezzi pubblici in Lombardia	50
4.1	La matrice OD	50
4.1.1	La matrice OD Lombarda	51
4.2	Analisi della matrice	53
4.3	I collegamenti dei mezzi pubblici Lombardi	56
5	Modello stocastico <i>SEII_aR</i>	58
5.1	Introduzione al modello	58
5.2	La costruzione del modello	59
5.2.1	Le equazione stocastiche	59
5.2.2	Mobilità e dinamica globale delle infezioni	62
5.3	La scelta dei parametri	63
5.4	Il modello applicato alla rete Lombarda	64
	Conclusioni	66
	Appendice	68
	Bibliografia	91

Elenco delle figure

1.1	Curve epidemiche Italia e Hubei, in rosso i picchi	11
1.2	Curva incidenza Hubei, in rosso il suo picco	12
1.3	Curva incidenza Italia, in rosso il suo picco	13
2.1	Intervallo Seriale Sars vs Covid-19	25
3.1	Modello SIR	37
3.2	SIR Lombardia	40
3.3	Differenze adattamento modelli ai dati	42
3.4	Andamento effettivo R_t in Lombardia	44
3.5	Modello SEIR	45
3.6	SEIR Lombardia	48
4.1	Mezzi di trasporto più utilizzati in Lombardia	54
4.2	Motivi di spostamento in Lombardia	55
4.3	Mezzi di trasporto più utilizzati in Lombardia per motivo	55
4.4	La rete dei trasporti dei mezzi pubblici Lombardi	56
5.1	Funzionamento modello stocastico di tipo $SEII_aR$	60
5.2	Andamento incidenza tramite modello stocastico	64

Introduzione

Malattie ed epidemie da sempre sono una delle maggiori cause di morte nel mondo, anche più di guerre e carestie. Gli agenti patogeni si adattano e si evolvono nel corso del tempo, determinando così, la comparsa di nuove malattie o il rigenerarsi di antiche forme di contagio che si pensavano estinte. È chiaro che l'invasione umana con nuovi ecosistemi, fattori quali il surriscaldamento globale, il degrado ambientale, il potenziamento degli spostamenti e comunicazioni su scala internazionale, e i cambiamenti in campo economico, abbiano cominciato a fornire sempre più terreno fertile per la generazione di nuove malattie.

Dati gli eventi accorsi in questo inizio anno, si è deciso di destinare il lavoro di tesi sullo studio del *Covid-19*, meglio conosciuto come coronavirus. La malattia infettiva, causata da un virus mai identificato prima negli esseri umani, ha prodotto sconcerto, morti e misure straordinarie (atte al contenimento e al rallentamento della diffusione) in ogni parte del mondo come mai prima ad ora nel dopo guerra. L'idea di questo lavoro nasce prima dell'arrivo della malattia nel territorio italiano quando ancora i soli casi erano limitati al territorio cinese. Verranno affrontati capitolo per capitolo i passi necessari per la costruzione di un modello epidemiologico in grado di stimare come si sarebbero trasmesse e con che velocità le infezioni in Lombardia sulla base degli spostamenti effettuati dalle persone tramite mezzi pubblici. Lo scopo finale è di capire come i decreti abbiano limitato e cambiato la diffusione della malattia cercando di capirne il mese di arrivo nel suolo Lombardo.

Nel capitolo 1 si affronterà l'analisi dei dati. Il primo passo sarà dato dalla costruzione della curva epidemica del virus, fondamentale per tutte le analisi successive, per poi affrontare le varie misure utilizzate nel mondo per l'analisi dei dati attraverso una loro visione critica. Come ultimo passo in questo capitolo si troverà il rischio di morte della malattia andandolo a confrontare con quello dell'influenza, mettendo chiarezza a quanto mal comunicato da media e giornali in questi ultimi mesi.

Il capitolo 2 costituisce una parte importante del lavoro. Si affronterà la spiegazione e la delucidazione del parametro R_0 , con le relative differenze con R_t . Si stimerà tale parametro, fondamentale per capire la velocità di diffusione della malattia, tramite due tecniche statistiche: "*tasso di crescita esponenziale e massima verosimiglianza*".

Lo studio verrà effettuato in parallelo sia per la Cina che per l'Italia, per trovare e comparare eventuali differenze. Il parametro sarà inoltre necessario all'interno dei modelli.

Nel capitolo 3 saranno costruiti i modelli quali SIR e SEIR per cercare di comprendere come si sarebbe diffusa l'epidemia senza l'intervento dei decreti emessi. Rielaborando le equazioni differenziali del modello SIR e integrandole a tecniche quali regressione lineare e di Poisson si stimerà e si valuterà l'andamento nel tempo del parametro R_t .

Nel capitolo 4 verrà affrontato lo studio di una particolare matrice il cui nome è *Origine-Destinazione*, contenente informazioni relative agli spostamenti degli individui durante un giorno lavorativo. L'analisi sarà effettuata sulla Lombardia, con una breve analisi generale, per l'estrapolazione di alcune informazioni di interesse. Si costruirà poi, tramite le informazioni presenti all'interno della matrice, la rete dei trasporti pubblici nel territorio Lombardo.

Nell'ultimo capitolo si applicherà un modello epidemiologico complesso di tipo stocastico che riflettendo gli spostamenti degli individui in Lombardia cercherà di valutare un possibile impatto del virus senza la considerazione dei vari decreti emessi dal governo (hanno limitato il numero degli spostamenti). Tramite questo modello si cercherà di capire quando il Covid-19 è effettivamente arrivato in Italia.

Si riporta nell'appendice il codice utilizzato nell'ambiente di lavoro RStudio.

Capitolo 1

Analisi dei dati

1.1 L'analisi descrittiva dei dati

La descrizione accurata dell'epidemia costituisce il primo importante passo per comprendere un virus. Un'analisi descrittiva dell'episodio può bastare per indirizzare le misure immediate di controllo a sviluppare ipotesi sulla sorgente di infezione e sulle modalità di trasmissione. Per questi motivi il primo passo è quello della costruzione di misure, come quella della curva epidemica (uno tra i più utilizzati), che diano informazione riguardanti la situazione reale.

Ponendo in un grafico il numero di casi (incidenza) e il tempo si osserverà una curva, detta epidemica, nella quale sarà possibile osservare la crescita dell'infezione[1]. Lo scorrere del tempo è visualizzato sull'asse delle x mentre i conteggi dei casi vengono visualizzati sull'asse delle y . Il risultato è una rappresentazione visiva dell'insorgenza della malattia nei casi associati a un'epidemia. È utile perché può fornire molte informazioni tra cui:

- Descrizione della diffusione
- Dimensione dei contagi nel tempo
- Valori anomali
- Periodo di incubazione della malattia

Tali aspetti dal punto di vista matematico/statistico possono essere spiegati in base al tipo di analisi che viene effettuata.

La descrizione della diffusione può essere osservata tramite misure come la curtosi, che rappresenta un allontanamento dalla normalità¹ distributiva, rispetto alla quale si verifica un maggiore appiattimento o un maggiore allungamento. In base all'appiattimento della curva si possono comprendere fenomeni come la velocità della diffusione della malattia. Una distribuzione platicurtica è indice di diffusione lenta. Una distribuzione leptocurtica rappresenta grossi problemi, in quanto essendoci diffusioni rapide, il sistema sanitario di un determinato paese potrebbe non farcela.

La dimensione dei contagi nel tempo di una malattia può essere trovata visualizzando i picchi² di una determinata curva epidemica. Esistono due tipi di picchi per ogni indicatore, che viene utilizzato. Definendo un individuo infetto come colui che è contaminato da microrganismi infettivi, e che può quindi trasmettere o provocare un'infezione, si possono trovare i relativi picchi di questo tipo di misura definendo $f(t)$ come una funzione che descrive il numero dei casi infetti al tempo t :

- 1 Il picco dei nuovi casi (picco dei casi di incidenza): coincide con il conteggio più alto dei nuovi casi osservati di una malattia in un determinato istante temporale (ad esempio un giorno se ci si riferisce a dati analizzati giornalmente)
- 2 Il picco epidemico rappresenta il punto di massimo dei casi di prevalenza, ovvero il più alto numero di casi totali presenti in una popolazione di interesse, osservati durante l'epidemia stessa

L'individuazione dei valori anomali assume alta importanza nell'utilizzo delle varie misure. Un cambio di metodologia nella rilevazione degli individui infetti, ad esempio aumentando il numero di tamponi effettuati ogni giorno, può portare a interpretazioni sbagliate. Se nel giorno t vengono effettuati 2000 tamponi (si supponga che nei giorni $t - 1$, $t - 2$... siano stati sempre lo stesso numero), di cui 1000 positivi, e nel giorno successivo $t + 1$ ne vengono effettuati 9000, di cui 2000 non negativi, il confronto diretto di persone trovate positive non avrebbe alcuna rilevanza statistica senza tenere in considerazione il numero di tamponi effettuati. Sarebbe infatti sbagliato affermato che nel periodo $t + 1$ siano presenti più casi rispetto ai precedenti. È quindi importante nel caso vengano rilevate osservazioni anomale di cercare di comprendere il motivo che le rende tali (nell'esempio significherebbe indagare il motivo dei più tamponi effettuati) per valutarne una loro esclusione dall'analisi.

Nell'incidenza, che rappresenta i nuovi casi giornalieri e che viene rappresentata da istogrammi (grafico a barre), la distanza temporale tra i vari massimi relativi (barre più alte) è molte volte indice del periodo di incubazione del virus.

¹una curva epidemica è approssimabile a una distribuzione normale [2]

²i differenti picchi verranno esposti in maniera più chiara tramite visualizzazione grafica nei successivi due paragrafi

1.2 La prevalenza: Italia e Hubei

Un tipo di curva epidemica, che rappresenta uno degli strumenti più utilizzati insieme all'incidenza, è quella data dalla prevalenza (indicatore dei nuovi casi giornalieri). La prevalenza indica il numero di individui (a volte anche espresso in termini percentuali) di una certa popolazione che riscontra una malattia in un determinato momento temporale. Se si verificano nuovi casi di persone infette la curva avrà tendenza crescente. Tuttavia se gli individui con la malattia muoiono (il tasso di mortalità è differente da zero) la prevalenza diminuirà di conseguenza. Anche il numero dei guariti incide su questo calcolo nella stessa maniera del numero delle morti e quindi questo potrebbe costituirne un limite. Pertanto è importante esaminare oltre alla prevalenza altre tipo di curve come l'incidenza.

È bene chiarire subito i due concetti con un esempio. Se si pensa ad una vasca da bagno l'incidenza è l'acqua che viene aggiunta nella vasca ogni lasso di tempo t mentre la prevalenza è il contenuto della vasca all'istante di tempo t e i casi che muoiono/guariscono sono rappresentati dall'acqua che esce per evaporazione.

Si riassume quanto detto in termini algebrici elementari:

$$X(t) = N(t) - M(t) - G(t) \tag{1.1}$$

Ove:

- $X(t)$ sono le persone attualmente positive al tempo t
- $N(t)$ il numero totale di casi al tempo t
- $M(t)$ il totale dei morti al tempo t
- $G(t)$ il totale dei guariti al tempo t

La curva epidemica basata sulla prevalenza per l'Italia e Hubei³ è data dalla figura 1.1. Sull'asse delle y si trovano gli attualmente positivi mentre sull'asse delle x si ha l'andamento temporale.

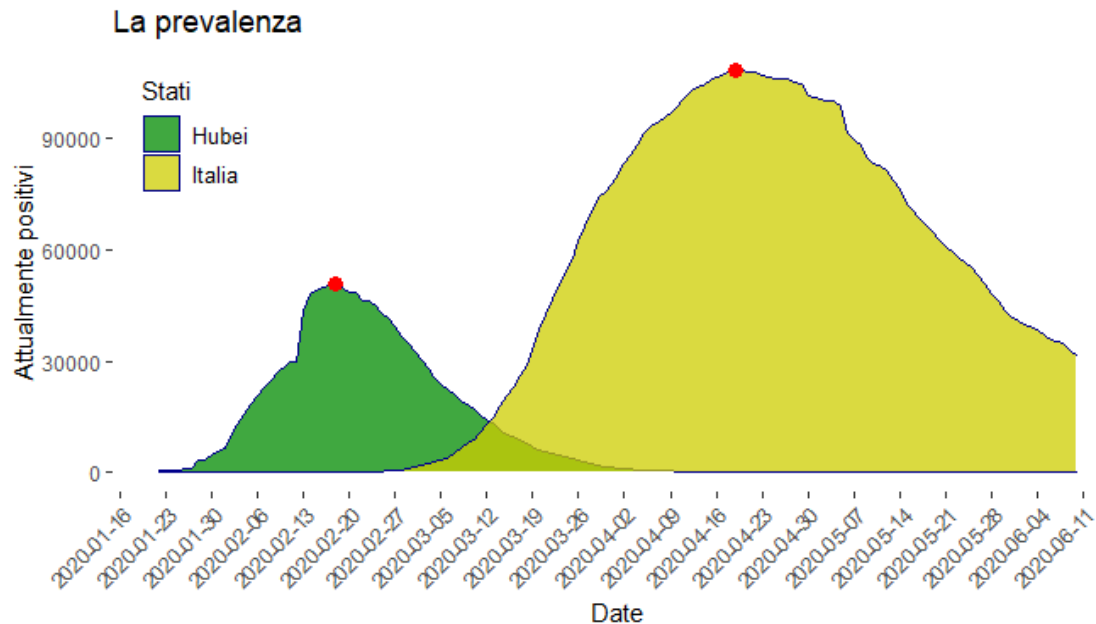


Figura 1.1: Curve epidemiche Italia e Hubei, in rosso i picchi

I periodi di inizio del virus sono differenti così come la loro diffusione. Il picco più alto, indice delle persone attualmente positive in un certo istante di tempo (giorno), indica che la dimensione del contagio è stata maggiormente sviluppata nel territorio. In parole semplici, in Italia si è visto un raggiungimento di positivi maggiore rispetto ad Hubei. La curva cinese appare inoltre meno piatta, il che vuol dire che, come spiegato in precedenza, la diffusione è stata più veloce ed è terminata prima.

³provincia Cina sede del primo focolaio

1.3 Incidenza: Italia e Hubei

La prevalenza differisce quindi dall'incidenza della malattia che rappresenta il numero di nuovi casi che si sviluppano in un determinato periodo di tempo⁴. L'incidenza permette di verificare l'andamento nel tempo giornaliero dei nuovi casi all'interno del territorio di riferimento.

Si parte a visionare questa misura per quanto riguarda la provincia cinese. In questo caso sull'asse delle x si ha sempre il tempo mentre sulle y si hanno i valori dei nuovi casi positivi giornalieri. Si aggiunge anche un grafico con finestra mobile pari a 3 giorni per aver una miglior visuale del trend.

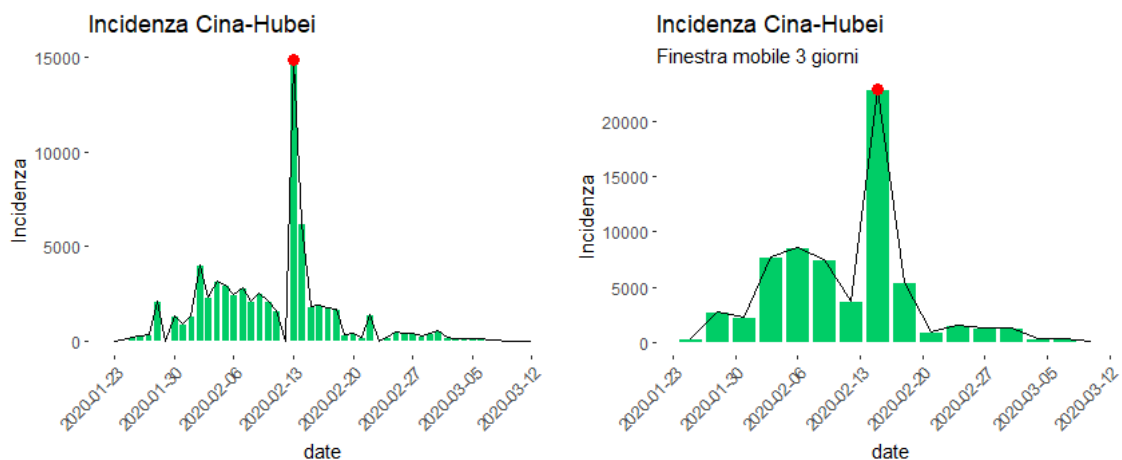


Figura 1.2: Curva incidenza Hubei, in rosso il suo picco

Il picco di nuovi casi facilmente osservabile, rappresenta un valore anomalo che si spiega con un cambio di metodologia di rilevazione dei positivi. Infatti in quel giorno i tamponi sono stato effettuati su larga scala a differenza del giorno precedente, dove i tamponi erano riservati solo ai casi sospetti di individui con sintomi. La considerazione è importante in quanto il cambio di rilevazione comporta errori nell'analisi se non se ne tiene conto. L'andamento dei nuovi contagi è quasi costante nel tempo. Non si osserva una particolare crescita esponenziale, molto probabilmente non rilevata dal metodo di tamponamento limitato e decresce gradualmente fino ad azzerarsi (indice di positivi pari a zero).

⁴il tempo utilizzato per il calcolo dei nuovi casi è di un giorno

Visionando l'incidenza nel territorio italiano si trova invece un andamento differente. Si aggiunge anche in questo caso una finestra mobile di tre giorni per visualizzare al meglio il trend. Il colore nel grafico rappresenta il numero dei tamponi effettuati giornalmente/ogni tre giorni rispettivamente nel grafico senza finestra mobile e con finestra mobile.

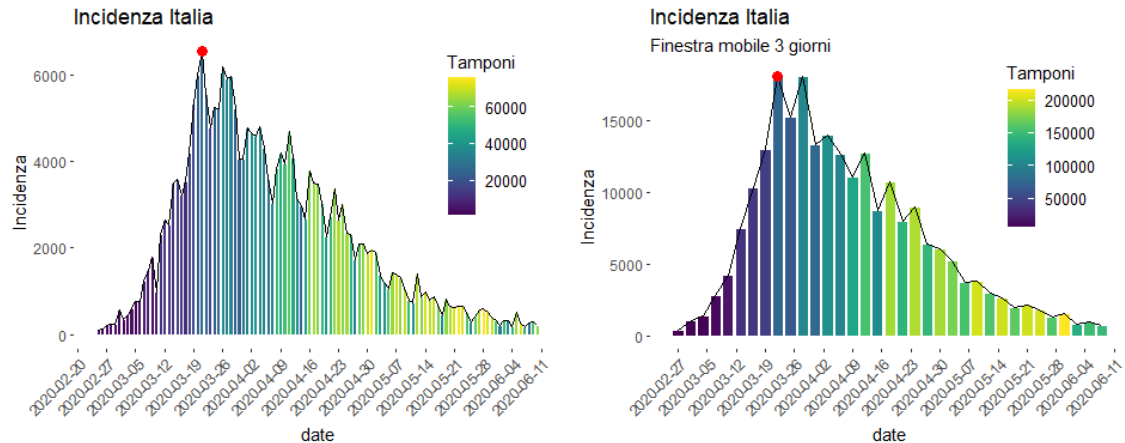


Figura 1.3: Curva incidenza Italia, in rosso il suo picco

Dalla curva si può notare un andamento esponenziale di crescita molto più omogeneo rispetto a quello di Hubei. Si vede inoltre che la crescita sembra arrestarsi intorno al primo aprile, a seguito probabilmente delle misure attuate dal governo. Quello che appare molto interessante è che i vari picchi si trovano in corrispondenza di un intervallo temporale di circa 5-6-7 giorni che corrisponde al periodo medio di incubazione del Covid-19.

Talvolta una variazione dell'incidenza può essere soltanto apparente ed essere dovuta a un artefatto come nel caso in cui vengano utilizzati nuovi metodi diagnostici che consentono di apprendere la malattia già al suo stato precoce. In questo caso anche se l'incidenza reale della malattia rimane costante, attraverso il nuovo strumento diagnostico verrà individuato un numero maggiore di casi rispetto a quanto si era verificato in passato e ciò corrisponderà appunto a un aumento dell'incidenza soltanto apparente. In questo caso si ha un andamento costante dal 15 al 30 marzo circa a fronte però di un numero di tamponi ogni giorno più elevati, il che indica che l'incidenza in realtà non è costante ma in diminuzione in tale periodo.

Un'ultima osservazione ricade sulla parte finale del grafico in cui si nota che a fronte di un numero sempre più elevato di tamponi l'incidenza diminuisce.

1.4 Analisi delle varie misure e problematiche annesse

Ogni qual volta si fa affidamento ad un certo tipo di misura bisogna capire ed analizzare che tipo di metodo si sta utilizzando per comprenderne i relativi limiti. Si riassumono le misure più utilizzate in ambito dello studio del Covid-19 (utilizzate quotidianamente per la comunicazione dei dati), le loro caratteristiche e i loro limiti [3]:

- **Frequenza cumulata di nuovi casi:** È l'indicatore maggiormente diffuso e utilizzato. Definendo la frequenza cumulata, associata ad una modalità⁵ (in questo specifico caso trattasi di nuovi casi giornalieri), come la somma della sua frequenza assoluta e di quelle delle modalità che la precedono, questo indicatore rappresenta un andamento cumulativo di casi giornalieri. Questo metodo è molto debole dal punto di vista rappresentativo perchè dipende dal numero dei tamponi che vengono effettuati quotidianamente. Inoltre il numero dei nuovi casi del giorno t rappresenta in realtà i casi del giorno $t - 1$ in quanto i risultati dei tamponi sono disponibili soltanto il giorno seguente[4]
- **Il numero dei decessi giornalieri:** La qualità del dato è pessima. I risultati di questi numeri comunicati quotidianamente da protezione civile si riferiscono solamente ai casi di morte certa con Covid-19 all'interno degli ospedali. Non vengono quindi considerate tutte le morti delle persone all'interno delle abitazioni, che non sono poche, e le morti dovute al coronavirus non certificate dal tampone positivo a causa del loro numero limitato
- **Il numero totale delle morti:** Tramite frequenza cumulata è impossibile individuare un trend crescente/decescente dei deceduti, visibile soltanto tramite il calcolo di morti giornaliere che per i motivi descritti in precedenza non sono attendibili. Il numero totale delle morti è quindi poco utile a fornire un'immagine precisa dell'andamento della pericolosità della malattia.
- **Numeri registrati su una scala logaritmica:** È una misura che visualizza i dati trasformati in scala logaritmica su un piano cartesiano. L'asse delle y viene etichettato ad esempio da 1, 10, 100, 1000 mentre l'asse delle x rappresenta il tempo. Questo tipo di misura è utile per la visione e il confronto dei vari trend che esprimono l'andamento degli infetti tra paesi e nazioni ma non viene utilizzato, in quanto poco significativo, per valutare l'andamento della malattia

⁵in statistica una modalità di un carattere è uno dei suoi possibili valori

- **Previsione tramite modelli:** Modelli statistici più o meno avanzati vengono utilizzati per prevedere sia l'andamento dei casi infetti (per capirne i giorni necessari per il raggiungimento del massimo assoluto all'interno della curva epidemica) sia per capire entro quanto tempo i contagi raggiungono soglia zero. Costituiscono una grande importanza all'interno dello studio di una malattia. Una previsione rende consapevoli i paesi dei vari rischi e da quei risultati, capendo la portata della malattia, possono essere prese varie decisioni come la quarantena, chiusura dei mezzi pubblici....
Il problema di tali modelli è che sono molto imprecisi il che può essere spiegato dall'incertezza dei dati raccolti. È molto difficile prevedere gli infetti futuri sulla base di un numero di infetti, raccolto tramite tamponi, che è per certo sbagliato. Ad oggi infatti nessun modello è stato in grado di predire correttamente quello che poi è successo.
Alcune volte comunque vengono riportati dai media e telegiornali modelli senza incertezze sulle previsioni e senza spiegarne le problematiche annesse
- **Letalità del virus:** Spesso vengono citate percentuali di letalità del virus generiche che non considerano l'età, i problemi degli individui non sani e il genere dell'individuo, il che costituisce un'informazione non completa. Non si conosce il motivo ma i dati parlano in modo chiaro, apparentemente muoiono più maschi che femmine (*donne* $\approx 41\%$), e l'età colpita maggiormente è quella degli over 65 con problematiche di salute già esistenti. L'età media dei pazienti deceduti e positivi in Italia, i cui dati sono però molto simili nel mondo, è di 80 anni. È stato inoltre trovato su un campione di 3200 deceduti che soltanto il 4% di loro non aveva patologie pregresse[5]

1.5 La mortalità in Italia

In Italia in data 29/05/2020 si contano 33 229 decessi. Media, telegiornali e politici hanno e comunicano ancora oggi tantissimi numeri riguardanti il rischio di mortalità, tutti sempre differenti fra di loro, portando a molta confusione. Bisogna quindi capire il motivo di questi risultati differenti anche se è evidente che in base a come viene calcolata la mortalità, i risultati saranno completamente diversi.

Si parte da un esempio semplice per poi fare le stesse considerazioni nel territorio Italiano. Per capire meglio il concetto si supponga che si riesca a sapere in qualunque istante di tempo il numero esatto di persone infette, sia con sintomi sia senza (le persone che non presentano sintomi si definiscono asintomatiche). La popolazione di riferimento puramente immaginaria si suppone essere composta da 100 persone all'istante t_0 . All'istante t_1 si scopre che sono presenti 20 persone infette di cui 7 hanno presentato i sintomi della malattia e 13 no. Si considera infine un tempo t_2

dove si hanno ancora lo stesso numero di persone infette del tempo t_1 con la differenza che tra quelle persone sintomatiche, 5 sono guarite (una persona si definisce guarita quando dopo la presenza dei sintomi ritorna in stato di salute) e 2 sono morte.

In questo esempio si calcola il rischio di morte con i metodi standard che sono i seguenti:

- **Case fatality rate (CFR):** È la percentuale di decessi per una determinata malattia rispetto al numero totale di persone che l'hanno riscontrata e ne hanno manifestato i sintomi riuscendo a guarire. Non vengono quindi considerati in questo calcolo gli asintomatici. Si intuisce quindi come il numero finale sia il più delle volte alto. Questo tipo di indice spesso viene confuso con il rischio di morte ma in realtà è usato soltanto per avere informazioni sui deceduti e/o sul funzionamento delle varie misure restrittive adottate. E' anche utile per capire in quali fasce di età si rilevano un maggior numero di decessi.

Si può definire il case fatality rate sulla base dell'esempio precedente in questo modo:

$$\text{Case fatality rate} = \frac{\text{morti}_{t=2}(2)}{\text{guariti}_{t=2}(5)} = 0.4 = 40\%$$

- **Infection fatality rate (IFR):** A differenza del precedente è un indice di rischio di morte ma solo se si è infetti. Rappresenta semplicemente il rapporto tra morti e persone risultate positive al virus, in questo caso quindi vengono inclusi anche gli asintomatici.

In termini algebrici molto elementari si può calcolare per l'esempio precedente:

$$\text{Infection fatality rate} = \frac{\text{morti}_{t=2}(2)}{\text{infetti}_{t=2}(20)} = 0.10 = 10\%$$

In questo caso quindi si ha che una persona infetta ha la probabilità del 10% di morire

- **Crude mortality rate (CMR):** Come il precedente è un indice di rischio di morte ma al posto di essere confrontato con tutti gli infetti in un certo istante di tempo viene confrontato con l'intera popolazione presente all'istante di tempo iniziale.

$$\text{Crude mortality rate} = \frac{\text{morti}_{t=2}(2)}{\text{popolazione}_{t=0}(100)} = 0.02 = 2\%$$

Quest'ultimo è l'indice corretto per capire quanto sia alta la letalità di una malattia.

Ora si procede mostrando come i vari indici sono utilizzati all'interno del contesto del territorio italiano. In primo luogo si procede col mostrare l'utilità del **case fatality rate**, cercando di capire se i decreti hanno effettivamente abbassato questo tasso. I tempi considerati sono istante temporale $t = 10/03/2020$ e $t = 10/05/2020$. I risultati sono i seguenti:

$$CFR_{t=10/03/2020} = \frac{morti = 631}{guariti = 1004} = 0.628 \approx 63\%$$

$$CFR_{t=10/05/2020} = \frac{morti = 30560}{guariti = 105186} = 0.290 \approx 29\%$$

Il tasso quasi dimezzato indica come i decreti abbiano aiutato alla riduzione delle morti, ma è totalmente sbagliato definire come rischio di morte uno dei due risultati. Si procede ora con un confronto diretto tra **infection fatality rate** e **crude mortality rate** entrambi calcolati all'istante temporale $t = 29/05/2020$. Di solito è sempre preferibile il calcolo di questi indici quando la malattia cessa di esistere, perchè rende il risultato *definitivo*. Calcolandolo invece con intervalli temporali si ha un valore sicuramente differente (anche se molto simile), che però permette di avere un'idea molto precisa.

$$IFR_{t=29/03/2020} = \frac{morti = 33229}{infetti = 152844} = 0.22 \approx 22\%$$

$$CMR_{t=10/05/2020} = \frac{morti = 33229}{popolazione = 60,36mln} = 0.00056 \approx 0.056\%$$

Il risultato dato dall'infection fatality rate è stato spesso usato dai telegiornali e può creare preoccupazioni, ci dice infatti che una persona infetta ha il 22% di probabilità di morire (chiaramente in quel periodo di tempo). Bisogna però evidenziare il limite di questo calcolo che nell'esempio effettuato in precedenza (quello della popolazione di 100 persone) è stato sorvolato ipotizzando di conoscere con certezza il numero di infetti in ogni istante temporale. Questo come ben noto non è possibile, essendo i tamponi giornalieri limitati, e quindi il suo risultato è privo di significato. Non ha senso utilizzarlo ne per capire se le politiche adottate sono corrette ne per confrontare i rischi di morte con altri paesi, è facile intuire infatti come ogni paese abbia un numero di tamponi differente.

Per questo motivo quello maggiormente utilizzato è il CMR che può essere confrontato anche con differenti malattie per capire le loro differenze. Un confronto doveroso è quello con l'influenza per capirne se effettivamente il rischio di morte è maggiore. Il CMR dell'influenza è il seguente:

$$CMR_{t=2019} = \frac{morti \approx 10000}{popolazione = 60,36mln} = 0.00016 \approx 0.016\%$$

Il risultato che ne deriva è che il Covid-19 ha una letalità di 0.035 punti percentuali maggiore.

Come ultima considerazione bisogna osservare che il numero delle morti di coronavirus in Italia pari a 33 229, ma così come nel resto del mondo, è un numero basato sui morti che abbiano effettuato un tampone risultato positivo al virus. A molti deceduti non è stato effettuato quest'ultimo, forse per mancanza di strumenti a sufficienza, ma quello che ne risulta è un numero comunque *sottostimato*. Bisogna inoltre considerare le persone che sono venute a mancare direttamente nelle loro abitazioni che chiaramente non hanno effettuato il test. È quindi formalmente più corretto definire che il coronavirus ha una letalità ≈ 3.5 volte maggiore di quella dell'influenza, che potrebbe essere leggermente più alta.

Capitolo 2

La stima di R_0

2.1 Introduzione ai modelli epidemiologici

I modelli svolgono un ruolo importante nella comprensione e gestione delle dinamiche di trasmissione dei vari agenti patogeni. Possono essere utilizzati per descrivere nel tempo l'evoluzione della malattia, nonché per esplorare e/o comprendere meglio i fattori che influenzano l'incidenza dell'infezione. La modellizzazione è quindi un passo fondamentale per comprendere quali trattamenti e interventi possono essere più efficaci e quali fattori specifici devono essere considerati quando si cerca di combattere un'infezione.

Per comprendere le dinamiche complesse alla base della trasmissione delle malattie, si usano spesso una serie di modelli chiamati *modelli compartimentali*. I modelli compartimentali (di cui ci si occuperà in questo lavoro) sono una tecnica utilizzata per semplificare la modellazione matematica di malattie infettive. La popolazione viene divisa in compartimenti con il presupposto che ogni individuo nello stesso scomparto abbia le stesse caratteristiche[6]. Questi tipi di modelli sono di solito formati da equazioni differenziali (che sono deterministiche), ma possono anche avere una forma stocastica¹, che è una rappresentazione più realistica ma più complicata da analizzare. In questo caso i modelli deterministici modellano la diffusione tramite parametri fissi mentre quelli stocastici introducono variabili casuali (come la binomiale) per spiegare il "successo" o "insuccesso" di un'infezione.

I modelli compartimentali possono essere utilizzati per prevedere la grandezza di diffusione della malattia oppure per prevederne il raggiungimento del picco.

¹modello matematico adatto a studiare l'andamento dei fenomeni che seguono leggi casuali, probabilistiche

Bisogna comunque fare molta attenzione alle caratteristiche individuali delle persone che talvolta sono differenti, ad esempio un individuo può essere asintomatico, può contrarre più facilmente una malattia perchè più anziano...

Si è quindi detto che il funzionamento dei modelli avviene secondo una precisa suddivisione della popolazione. Qui di seguito si riportano le categorie in cui le persone, appartenenti ad una certa popolazione, vengono suddivise:

- **Suscettibili:** Individui che possono essere infettati
- **Esposti:** Individui infettati che a causa del periodo di incubazione dell'agente patogeno, non sono ancora infettivi e non hanno manifestano sintomi
- **Infetti:** Individui infetti da un agente patogeno ed in grado di trasmettere l'infezione ad altri. Un individuo è infetto quando presenta i sintomi
- **Guariti:** L'individuo non è più contagioso o è "rimosso" dalla popolazione (con rimosso si intende la sua morte)

2.2 Il parametro R_0

Si è spiegato in maniera semplicistica, il funzionamento dei vari tipi di modelli che chiaramente necessitano di alcuni parametri che si definiscono "fondamentali" senza i quali non si possono inizializzare i vari algoritmi. Naturalmente tanto più complesso è un modello tanto più i parametri diventano numerosi e difficili da calcolare. Uno di questi però indipendentemente dal modello utilizzato compare sempre e si definisce **basic reproduction number**, più conosciuto come R_0 .

R_0 è un termine matematico che rappresenta una misura della velocità con cui una malattia infettiva progredisce inizialmente^[7]. R_0 è caratteristico dell'epidemia nella sua fase iniziale, quindi per capire se le misure di contenimento hanno effetto, viene considerato il numero di riproduzione effettivo, indicato con R_t ², in quanto si riferisce ad un giorno generico t . R_0 è sicuramente importante, perché fornisce l'informazione sintetica di quanti casi secondari (casi che si verificano in tempi successivi al periodo di infezione del caso primario³ che derivano dal contatto con essi) vengono generati. Risulta infatti ovvio che quando un'infezione si diffonde a nuove persone, si riproduce. Il parametro fornisce quindi un'informazione sintetica di

²si spiegherà poco più avanti la loro relazione nel dettaglio

³la terminologia può applicarsi solo a malattie infettive che si diffondono da uomo a uomo, e si riferisce alla persona/gruppo che per prima/i trasmette/ono una malattia ad un gruppo di persone

quanti casi secondari vengono generati, per trasmissione interpersonale, da un caso primario in una popolazione completamente suscettibile. Ne consegue che una epidemia si instaura quando per ogni caso primario si generano più casi secondari e da ognuno di questi vengono generati altri casi. All'opposto, se ogni singolo caso non ne contagia nessun altro, la circolazione dell'infezione è destinata ad estinguersi. Intuitivamente è apprezzabile che il valore è direttamente proporzionale al numero di contatti per giorno del caso primario (più persone incontra, più persone si infettano), alla durata della sua fase di contagiosità (più a lungo rimane contagioso, più è alto il numero delle persone che contagia) e alla probabilità di trasmissione dell'infezione per singolo contatto (un individuo infetto non è detto che trasmetta la malattia a tutti gli individui che incontra). Tutte queste quantità sono difficili da osservare direttamente e in genere ci si basa su stime, sotto diverse assunzioni, che vengono utilizzate per la costruzione dei modelli statistici che si affronteranno nei prossimi capitoli. La spiegazione è che risulta impossibile avere a disposizione i dati sul numero di infetti giornalieri, per motivi come il numero di tamponi effettuati che risultano limitati. Se questa informazione fosse comunque nota allora il calcolo diretto sarebbe molto semplice, perchè basterebbe dividere gli individui infetti al periodo di tempo t_x con quello t_{x-1} (ove x arco temporale di uno dei primi giorni della diffusione della malattia).

Quando si osserva una trasmissione di contagi nella popolazione generale come il coronavirus, R_0 viene stimato in modo empirico, ossia osservando la velocità di crescita del numero totale dei casi giorno dopo giorno (stima tramite tasso di crescita esponenziale). Sapendo la data di insorgenza dei sintomi, il tempo di incubazione e l'intervallo di tempo tra la comparsa dei sintomi nel caso primario e la comparsa dei sintomi nei casi secondari (questo intervallo di tempo si definisce Generation-Time) è possibile ricostruire le diverse generazioni di casi e stimare l'indice di riproduzione.

Se una malattia ad esempio ha un R_0 di 1.30⁴, una persona che la contrae la trasmetterà ad una media di altre 1.30 persone. La lettura chiave che viene data a questo parametro è la seguente:

- $R_0 < 1$: Se R_0 è minore di 1 ogni infezione esistente provoca meno di una nuova infezione. In questo caso, la malattia diminuirà e alla fine si estinguerà
- $R_0 = 1$: Se R_0 è uguale a 1, ogni infezione esistente provoca una nuova infezione. La malattia rimarrà viva e stabile, ma non ci sarà un focolaio o un'epidemia
- $R_0 > 1$: Se R_0 è maggiore di 1, ogni infezione esistente provoca più di una nuova infezione. La malattia si diffonderà tra le persone e potrebbe esserci un focolaio o un'epidemia.

⁴ R_0 della comune influenza

Ripetendolo un'ulteriore volta, R_0 descrive la trasmissione in una popolazione completamente suscettibile. Ne deriva che per ottenere una stima bisogna studiare la malattia nella sua fase epidemica iniziale ove i contagi non vengono ridotti dalla quarantena. È importante sottolineare questa considerazione perchè i vari stimatori di R_0 possono essere utilizzati soltanto in questo caso, ovvero se nel campione di riferimento in cui si calcola l'indice:

- Nessuno è vaccinato
- Nessuno è immune alla malattia
- Non sono state prese misure di ordine pubblico per diminuire i contagi

Il motivo è presto spiegato, anche se una sola condizione non viene rispettata R_0 viene sottostimato perchè gli individui infetti non riuscirebbero a trasmettere il virus a parte della popolazione che in realtà, in mancanza di suddette ipotesi, si sarebbe infettata.

2.2.1 Alcune delucidazioni su R_0 e differenze con R_t

Leader mondiali e sanità pubblica hanno speso gli ultimi mesi e spenderanno i prossimi a parlare di questo R_0 . Può sembrare un concetto semplice ma in realtà è molto più complesso di quanto spiegato nel paragrafo precedente. È il risultato di metodi di stima molto complessi, e può cambiare radicalmente da un luogo ad un altro in base alle condizioni del territorio e ai comportamenti delle persone. Il parametro nasce da studi demografici dove è usato per descrivere il tasso di nascite⁵. R si riferisce alla riproduzione e 0 alla generazione 0esima (paziente zero).

Ripetendo l'esempio in precedenza se si suppone un valore del parametro $R_0 = 1.30$, se 1000 persone fossero infette ci si aspetterà di avere questo tipo di andamento di contagi nel tempo:

$$t_1 = 1300 \rightarrow t_2 = 1690 \rightarrow \dots \rightarrow t_{10} = 42621$$

Il metodo più utilizzato per stimare il parametro R_0 è una stima tramite tasso di crescita esponenziale ma ad oggi un consenso su come stimare il parametro non è presente. Per questo motivo c'è la necessità di associare oltre che ad una stima puntuale, un intervallo di confidenza (IC) in modo da avere una misura di incertezza sul valore. Molte delle volte infatti i vari intervalli di confidenza sono molto ampi. La

⁵si userà tale esempio per trovare l'equazione di Lokta-Eulero

SARS ad esempio, è di solito descritta con un \hat{R}_0 da 2 a 5. Il morbillo addirittura ha un IC che varia da 3 a 203, un enorme differenza. È evidente che gli scienziati stanno ancora discutendo e rivedendo in che modo possono essere trovate stime sempre più precise.

Il valore di R_0 può cambiare a seguito di modifiche nei contatti sociali (per esempio a seguito di interventi di distanziamento sociale) oppure a seguito della riduzione del numero di persone suscettibili. Viene allora solitamente indicato con R_t . Tale valore è la stima che viene comunicata giornalmente. La differenza è grossolana e non bisogna farne confusioni. Un andamento di R_0 nel tempo, che prende poi il nome di R_t serve ai governi per valutare l'effetto dei decreti che vengono emessi e può essere usato come parametro variabile all'interno di modelli statistici per previsioni future circa gli infetti. La valutazione della malattia viene espressa invece dalla sola stima del parametro R_0 che viene sempre effettuata solamente nella fase iniziale di una epidemia per i motivi spiegati in precedenza.

A tal proposito è bene sottolineare che se R_t è più piccolo di 1 non significa che il virus è sconfitto ma supponendo che i numeri siano corretti, che la sua diffusione è stata messa in pausa. In questo caso avendo un valore al di sotto della soglia, significa che ogni, per esempio, 100 malati si infetteranno meno di altri 100 individui. Ogni generazione successiva di infezioni sarà più piccola dell'ultima ma le persone potranno ancora ammalarsi e morire.

Per capire la rapidità della diffusione di una malattia, e per utilizzare gli stimatori di R_0 è necessario anche, oltre a questo parametro, la costruzione di un intervallo seriale (noto come *Generation-Time*) per costruire l'andamento delle diverse generazioni di casi nel tempo.

2.3 Generation-Time

Il tempo di generazione è definito come l'intervallo di tempo medio tra le infezioni di un caso primario e l'infezione di un caso secondario causata dal caso primario. La comprensione degli intervalli di tempo tra le generazioni successive tra individui infetti è fondamentale per quantificare in modo appropriato la dinamica di trasmissione delle malattie infettive.

Sebbene sia difficile osservare direttamente l'intervallo di tempo, in molti si occupano di derivarne andamenti il più precisi possibili per derivarne distribuzioni che giocano un ruolo chiave nella stima di trasmissione di una malattia (come spiegato

misurata da R_0). Il tempo di generazione è un periodo non osservabile⁶ che di solito viene chiamato anche con il termine "*intervallo seriale*", relativamente usati per distinguere la durata dall'infezione primaria a quella secondaria misurata in base ai tempi dell'infezione e all'insorgenza dei sintomi. Sia l'intervallo seriale che il tempo di generazione hanno però lo stesso valore medio e per questo ai fini di questo lavoro rappresentano la medesima definizione, sebbene come detto in qualche riga precedente hanno concezioni differenti.

In questo paragrafo si illustrano i risultati riguardanti le stime del generation-time per il coronavirus fornite da *International Society for infectious diseases*[8]. Il lavoro svolto per la ricerca della distribuzione si basa sull'osservazione di un campione di 28 coppie formate da infettore-infetto. La derivazione dei risultati ha implicato l'utilizzo di tecniche statistiche complesse che hanno messo a confronto tre tipologie di variabili casuali: *log-normale*, *gamma* e *Weibull*

Il risultato che ne è derivato è che il miglior modello per descrivere i vari intervalli di diffusione del virus è stato quello costruito con la log-normale con mediana dell'intervallo stimata attorno a 4.0 giorni con un intervallo di confidenza compreso tra [3.1, 4.9]. I parametri stimati della variabile casuale sono risultati:

$$\begin{aligned}\mu &= 4.7 \rightarrow IC[3.7, 6.0] \\ \sigma &= 2.9 \rightarrow IC[1.9, 4.9]\end{aligned}$$

Si ricorda che $f(x) \sim \log - normale(\mu, \sigma^2)$ ha funzione di densità:

$$f(x) = \frac{e^{-\frac{(\log x - \mu)^2}{2\sigma^2}}}{x\sqrt{2\pi}\sigma} \text{ per } x > 0 \quad (2.1)$$

La stima della mediana dell'intervallo seriale di 4,0 giorni indica che l'infezione da Covid-19 porta a rapidi cicli di trasmissione da una generazione di casi all'altra. Questo vuol dire che la sua diffusione è rapida e veloce. Se si confrontano i risultati con quelli della SARS (immagine 2.1 ove sull'asse delle x si ha l'andamento temporale e sull'asse delle y si ha la densità di probabilità) l'intervallo di trasmissione medio del coronavirus risulta molto più breve, il che è indice di una pericolosità di diffusione maggiore.

⁶poiché il momento preciso dell'infezione è molto difficile e quasi impossibile da rilevare, il tempo di generazione non è osservabile per due ospiti successivi

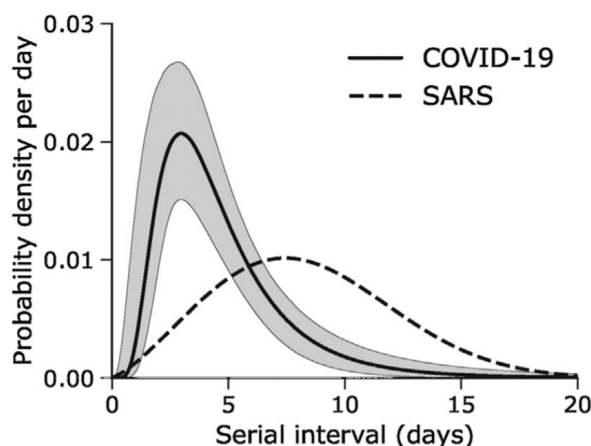


Figura 2.1: Intervallo Seriale Sars vs Covid-19

Una cosa ancora più importante è che la mediana dell'intervallo stimato è più breve delle stime sul periodo medio di incubazione (circa 5 giorni). Questo significa che è probabile che ci sia una trasmissione pre-sintomatica che può verificarsi più frequentemente della trasmissione sintomatica. Il risultato che ne è derivato in sostanza indica che molte trasmissioni non possono essere impediti solo mediante l'isolamento di casi sintomatici, poiché quando i contatti vengono rintracciati potrebbero già aver infettato e generato casi secondari e così via.

2.4 Stima di R_0 tramite il tasso di crescita esponenziale

2.4.1 Introduzione

Una stima accurata del valore di R_0 è cruciale per capire l'espansione di una malattia. Per le nuove malattie come lo stesso Covid-19, le informazioni disponibili sulla trasmissibilità sono purtroppo limitate al conteggio giornaliero dei nuovi casi. È noto come questi conteggi aumentano quasi in modo esponenziale nella fase iniziale di un'epidemia. Tra i metodi più utilizzati per il calcolo del parametro ci sono due correnti di pensiero differenti per la stima. La prima utilizza il tasso di crescita esponenziale, r , che è definito come il rapporto dei casi giornalieri nel tempo rispetto al giorno precedente. Il valore osservato di r viene quindi legato al valore R_0 attraverso una semplice equazione lineare[9]:

$$R_0 = 1 + rT_c$$

Ove T_c è il valore atteso della variabile casuale log-normale che definisce il generation-time (equazione 2.1).

Demografi, ecologi e biologi evoluzionisti invece adottano un approccio leggermente diverso, più matematico. Derivano R_0 secondo la cosiddetta equazione di **Lotka-Eulero**. In sostanza si semplifica questa equazione ignorando la variabilità del tempo di generazione ove per variabilità si intende la possibile modifica dei tempi di infezione fra due individui⁷ e anche lo spazio dei valori che l'intervallo di tempo medio può assumere sulla base dell'intervallo di confidenza.

Il risultato è un'equazione esponenziale:

$$R_0 = e^{rT_c}$$

Qui, T_c è un analogo demografico della media del generation-time epidemiologico.

Avendo due equazioni alternative per mettere in relazione il valore desiderato si affronta la difficoltà di scegliere quello più appropriato. Come è facile capire a seconda della scelta si avranno due risultati totalmente differenti.

Quello che viene fatto nella realtà è una riformulazione dell'equazione di Lotka-Eulero, che viene poi utilizzata da epidemiologi e statistici per la stima del parametro R_0 .

2.4.2 L'equazione di Lotka-Eulero

Si deriva l'equazione di Lotka-Eulero tramite degli esempi semplificativi usando la popolazione umana. Per semplicità, ci si concentra su individui di sesso femminile, supponendo che vi sia sempre una scorta sufficiente di maschi per garantirne la riproduzione. Si misura il tempo e l'età in anni e ci si riferirà al tempo presente come $t = 0$, in modo che gli eventi nel passato si siano verificati quando il tempo t è negativo e gli eventi nel futuro si verificheranno quando il tempo t è positivo. Si parte dal presupposto che la popolazione mostra una crescita esponenziale ad un tasso fisso e che la distribuzione per età della popolazione non cambi nel tempo.

L'equazione di Lotka-Eulero può essere intesa come la combinazione di due concetti che possono essere spiegati in modo molto intuitivo.

⁷è facile pensare in ambito di un virus che questo potendo mutare può variare le tempistiche di contagio tra individui

- 1- Se si sommano il numero di bambini nati da madri di tutte le età in un determinato momento, si ottiene il numero totale di nascite in quel momento. Il numero di nascite per madri di età a al tempo t è uguale al numero di nascite al momento $t - a$ (il numero di madri, comprese quelle che non sono sopravvissute) moltiplicato per il numero atteso di figli all'anno per le madri di età a . Sommando queste nascite in tutte le possibili età delle madri, si ottiene il numero totale di nascite nell'anno t ovvero:

$$b(t) = \int_{a=0}^{\infty} b(t-a)n(a) da \quad (2.2)$$

dove $b(t)$ si riferisce al tasso di natalità della popolazione al momento t e $n(a)$ si riferisce al tasso di produzione della prole femminile di una madre all'età a .

- 2- In secondo luogo, perché la popolazione sia in crescita esponenziale con una distribuzione per età stabile, il numero di nascite in un dato momento (ad esempio t) è uguale al numero delle nascite di un tempo fa, moltiplicato per la crescita esponenziale della popolazione. Si può quindi dedurre la seguente uguaglianza:

$$b(t) = b(t-a)e^{ra} \quad (2.3)$$

A questo punto si combinano le espressioni (2.3) e (2.4) per ottenere un'espressione con $b(t)$ su entrambi i lati:

$$b(t) = \int_{a=0}^{\infty} b(t)e^{-ra}n(a) da \quad (2.4)$$

Il parametro $n(a)$ è più familiare nella demografia come il prodotto delle funzioni di sopravvivenza al parto e fecondità, $n(a) = l(a)m(a)$. Usando questa parametrizzazione più familiare e rimuovendo $b(t)$ da entrambi i lati dell'equazione, si ottiene la cosiddetta equazione di *Lotka-Eulero*:

$$\text{Con quindi : } n(a) = l(a)m(a)$$

$$1 = \int_{a=0}^{\infty} e^{-ra}l(a)m(a) da \quad (2.5)$$

2.4.3 Derivazione di R_0 dall'equazione di Lotka-Eulero

Mentre l'equazione Lotka-Eulero è una parte fondamentale della demografia, con cui si potrebbe essere interessati a derivare i tassi di crescita della popolazione, il problema qui rilevante è stimare il numero riproduttivo, R_0 , dai tassi di crescita di una malattia. In precedenza si è definito $n(a)$ come il tasso di produzione della prole femminile da parte di una madre all'età a . Risulta che se si integra $n(a)$ per tutta la durata della vita di una madre, si ottiene il suo numero totale di figli. Tale numero altro non è che R_0 :

$$R_0 = \int_0^{\infty} n(a) da \quad (2.6)$$

Il tasso $n(a)$ può essere normalizzato a una distribuzione $g(a)$, che in questo caso rappresenta l'età in cui l'individuo è in gravidanza:

$$g(a) = \frac{n(a)}{\int_0^{\infty} n(a) da} = \frac{n(a)}{R_0} \quad (2.7)$$

Ora si considera l'età come il tempo trascorso per prendersi l'infezione, e quindi nella notazione sopra la distribuzione del *generation-time* è equivalente a $g(a)$. Sostituendo questa espressione con la distribuzione dell'intervallo seriale $g(a)$ nell'equazione di Lotka-Eulero (2.5), si ottiene saltando passaggi algebrici elementari:

$$\frac{1}{R_0} = \int_{a=0}^{\infty} e^{-ra} g(a) da \quad (2.8)$$

Il termine che ora appare nella parte destra di questa equazione è familiare ai matematici come la cosiddetta trasformata di Laplace della funzione $g(a)$. Più specificamente nota agli statistici come la funzione generatrice dei momenti $M(z)$ della distribuzione $g(a)$. Scritta formalmente è definita come:

$$M(z) = \int_{a=0}^{\infty} e^{zt} g(a) da \quad (2.9)$$

Si usa quindi la funzione generatrice dei momenti per semplificare la notazione dell'equazione riformata di Lotka-Eulero (2.8). Qui, l'argomento z prende il valore del tasso di crescita definito in precedenza ma ovviamente cambiato di segno ($-r$):

$$R_0 = \frac{1}{M(-r)} \quad (2.10)$$

Inutile dire che per trovare un valore di R_0 , $M(-r)$ deve esistere! Si ricorda che $M(z)$ funzione generatrice dei momenti essendo definita da una funzione integranda può accadere, infatti, che l'integrale sia divergente.

Una funzione generatrice di momenti (se esiste) caratterizza in modo univoco la forma dell'intera distribuzione di probabilità: $M(z)$ determina $g(a)$ e, al contrario, $g(a)$ determina $M(z)$. Il corollario di questo momento che genera l'espressione della funzione è quindi una relazione tra il tasso di crescita r e R_0 che caratterizza in modo univoco la forma della distribuzione del generation-time e, al contrario, la forma della distribuzione del generation-time determina la relazione tra il numero riproduttivo e il tasso di crescita.

Per le malattie con un intervallo medio T_c e una deviazione standard della distribuzione σ , come in questo caso, gli intervalli seriali possono approssimarsi ad una distribuzione normale. Supponendo un intervallo di generazione normalmente distribuito si ottiene la seguente relazione tra il tasso di crescita r e il numero riproduttivo R_0 :

$$R_0 = e^{rT_c - \frac{1}{2}r^2\sigma^2} \quad (2.11)$$

Questa relazione è una curva convessa che approssima una curva esponenziale. Si evince anche che una distribuzione più concentrata attorno alla media della distribuzione generation-time, cioè con un valore più basso per σ , porta a valori più alti per il basic reproduction number.

A questo punto l'ultimo step è quello di crearsi un intervallo di confidenza per il valore di R_0 per l'esigenza di associare alla stima puntuale una misura di incertezza. Quello che si fa è utilizzare una regressione di poisson sulla 2.10 per crearsi degli intervalli.

2.4.4 R_0 utilizzando campione di osservazioni cinesi e italiane

A questo punto ci si serve della curva epidemica ricavata nel capitolo 1, sia per quanto riguarda quello italiano sia quello cinese, e dell'intervallo seriale, per stimare il parametro per entrambe le nazioni.

La prima cosa da fare è scegliere un periodo temporale su cui ricavare le stime. Si è scelto di considerare in entrambi le nazioni i primi 15 giorni di propagazione del virus. La scelta non è stata casuale ma è dovuta al fatto che nei dati cinesi al giorno

20 circa si ha un picco inusuale di casi infetti dovuti all'inizio di tamponamenti del tutto casuali (prima era effettuato solo alle persone infette). Seppur si è pensato di considerare come valore anomalo tale osservazione, non avrebbe avuto comunque senso allungare l'intervallo in quanto le disposizioni poi emesse dal governo cinese hanno contribuito ad abbassare il numero di contagi e come ampiamente ripetuto c'è bisogno di un andamento esponenziale, senza restrizioni, per aver il risultato il più preciso possibile. Sulla base di queste premesse si è scelto, per poter fare poi effettuare un confronto, di considerare lo stesso periodo di tempo di 15 giorni anche per l'Italia.

In particolare per questa comparazione si sono considerati solo i dati di *Hubei*, provincia cinese che comprende Wuhan che è la provincia dove si è trasmesso maggiormente il virus. Questa scelta per avere una reale percezione dell'andamento dei contagi.

Si calcolano tramite (2.11) i valori di R_0 ottenendo i seguenti risultati:

$$R_0^{Cina} : 3.39[3.36, 3.42]^8$$

$$R_0^{Italia} : 3.60[3.53, 3.67]^9$$

Considerazioni:

- Sembrerebbe che l' R_0 italiano sia leggermente più alto, ma non si può affermare con certezza il risultato. Si deve però considerare che in Italia le misure attuate dal governo per contenere la diffusione sono state meno severe rispetto a quanto avvenuto in Cina e forse questa è una possibile causa del risultato ottenuto
- Si può affermare con certezza invece che la stima del parametro trovato di R_0 è in realtà sottostimato ma questo è facilmente intuibile. In sostanza i dati utilizzati (chiaramente anche da OMS¹⁰ e ricercatori) non sono veritieri perchè non comprendono tutte le persone che in un certo periodo hanno realmente riscontrato l'infezione ma soltanto quelle risultate positive ai tamponi. Non potendo fare un controllo totale della popolazione gli unici dati a disposizione sono questi. E' complicato anche lo studio di un sottocampione per poter poi applicare il risultato alla popolazione intera perchè presenterebbe molti problemi che derivano dalla non conoscenza del Covid-19. Senza precise informazioni uno studio del genere è molto complesso e difficile

⁸periodo di tempo considerato dal 23/01/2020 al 07/02/2020

⁹periodo di tempo considerato dal 22/02/2020 al 08/03/2020:

¹⁰organizzazione mondiale della sanità

- Ricordando che si è utilizzato come generation-time, una log-normale (equazione 2.1), con parametri fissi stimati si fa notare che facendoli variare secondo i valori degli intervalli di confidenza la stima di R_0 varia

Sebbene questo metodo di stima sia molto utilizzato per la sua facilità di solito si utilizza anche qualcosa di diverso. Il metodo è quello della funzione di massima verosimiglianza, che in ambito statistico è un approccio classico per la stima dei parametri.

2.5 Stima di R_0 tramite la funzione di massima verosimiglianza

2.5.1 Derivazione della stima di massima verosimiglianza

Si utilizza ora un metodo basato sulla verosimiglianza per stimare R_0 . Si userà ancora l'intervallo seriale (equazione 2.1) trovato in precedenza. Il metodo in questione prende il nome dai loro inventori White e Pagano[10]. Ci sono due tipi di massimizzazioni, la prima considera l'intervallo seriale noto, mentre il secondo si basa sulla non conoscenza dell'intervallo che viene stimato. Si userà il primo in quanto il generation-time è "noto" ma per la sua derivazione serve una spiegazione anche del metodo che stima entrambi.

Prima di illustrare la parte matematico statistica alla base dello stimatore si procede a elencare le sue ipotesi:

- Non vi siano casi importati da altri paesi nella popolazione di riferimento
- Non vi siano dati mancanti
- La popolazione si mescoli uniformemente
- La popolazione deve essere chiusa (no nascite e morti)
- Nei casi primari compaiono i sintomi prima dei loro casi secondari (solo per il metodo che stima anche l'intervallo seriale)

Per non complicare troppo le cose si supponga che le condizioni vengano rispettate. Come ben noto chiaramente non tutte le ipotesi sono vere nell'ambito italiano, mentre risultano in parte vere per la Cina se si considera solo il primo periodo di manifestazione del virus. In Italia infatti:

- Il pazienze zero sembra derivare dalla Germania
- Non è verificata l'attendibilità dei dati
- La popolazione è sotto vincoli governativi come la quarantena imposti dai primi giorni

I dati che utilizza lo stimatore non sono altro che i conteggi dei nuovi casi giornalieri (curva epidemica), definiti da un vettore $N = \{N_t\}$ e $t = 0, \dots, T$ con t che indica una certa unità di tempo e N_t il numero di nuovi casi all'istante di tempo t . Senza perdita di generalità, si assume che t indicizza il periodo temporale definito dai giorni¹¹.

Di norma all'inizio di un'epidemia l'intervallo seriale non è ben definito perchè non ci sono osservazioni a sufficienza per poterne ricavare stime. Il primo metodo proposto si basa quindi su una stima diretta sia del numero riproduttivo sia dell'intervallo seriale.

Come possibile modello, si ipotizzi che il numero di casi secondari prodotti da un individuo infetto segua una distribuzione di Poisson, con valore atteso R_0 , e che l'intervallo seriale sia descritto da una distribuzione multinomiale con vettore di probabilità p_1, \dots, p_k , con chiaramente $p_1 + \dots + p_k = 1$. In sostanza si cerca di modellare le infezioni da casi primari a secondari, da secondari a terziari e così via secondo delle probabilità che rappresentano il "successo" di infezione nei casi successivi. L'assunzione di una distribuzione di questo tipo implica che dopo un certo periodo di tempo, definito da un certo k , la probabilità di generare un caso secondario sia trascurabile.

All'interno del focolaio la situazione viene descritta nel seguente modo. N_0 sono gli individui che inizialmente manifestano l'epidemia. Ognuno di questi casi genera indipendentemente casi secondari secondo una distribuzione di Poisson con media R_0 . Si definisce poi X_0 come il numero totale di casi prodotti dai casi N_0 iniziali (il che significa $X_0 \sim Pois(N_0 R_0)$). All'interno del focolaio questi casi X_0 infettano altri individui nei k giorni successivi secondo una variabile casuale multinomiale (intervallo seriale definito in precedenza in questo modo).

¹¹tuttavia, questo metodo è applicabile a qualsiasi unità di tempo discreta

La notazione è quindi la seguente:

- N_i rappresenta il numero totale di casi il giorno i
- X_{ij} rappresenta il numero di casi presenti nel giorno j , che sono stati generati dai casi N_i
- X_i indica il numero totale di casi prodotti dai casi primari nel giorno i (ovvero $X_i = \sum_j X_{i,j}$)

Se k , la lunghezza massima dell'intervallo seriale, fosse per esempio tre, allora si può illustrare il tutto con il seguente schema:

$$\begin{array}{ccccccc}
 N_0 & & & & & & \\
 N_1 = & X_{01} & & & & & \\
 N_2 = & X_{02} & + & X_{12} & & & \\
 N_3 = & X_{03} & + & X_{13} & + & X_{23} & \\
 N_4 = & & & X_{14} & + & X_{24} & + & X_{34} \\
 N_5 = & & & & & X_{25} & + & X_{35} & + & X_{45} \\
 \vdots & & & & & & \vdots & & \vdots &
 \end{array}$$

Si noti che questi schemi non forniscono alcuna indicazione del momento in cui si è verificato il passaggio dell'infezione tra l'individuo primario e quello secondario, ma rappresentano solo il momento in cui i casi diventano sintomatici. È da notare che non viene osservato X_{ij} . Se si fosse a conoscenza di tale informazione si potrebbe facilmente stimare R_0 e il vettore di probabilità \mathbf{p} ($p = \{p_1, p_2, \dots, p_k\}$ che rappresenta la distribuzione dell'intervallo seriale), tramite massima verosimiglianza. La sua costruzione per la derivazione delle stime risulterebbe questa:

$$\begin{aligned}
 L(R_0, \mathbf{p} | \mathbf{N}, \mathbf{X}) = & \left[\frac{e^{-N_0 R_0} (N_0 R_0)^{X_0}}{X_0!} \right] \left[\left(\begin{array}{c} X_0 \\ X_{01} \cdots X_{0,1+k} \end{array} \right) p_1^{X_{01}} \cdots p_k^{X_{0,k}} \right] \times \\
 & \left[\frac{e^{-N_1 R_0} (N_1 R_0)^{X_1}}{X_1!} \right] \left[\left(\begin{array}{c} X_1 \\ X_{12} \cdots X_{1,1+k} \end{array} \right) p_1^{X_{12}} p_k^{X_{1,1+k}} \right] \cdots \\
 & \times \left[\frac{e^{-N_T R_0} (N_T R_0)^{X_T}}{X_T!} \right] \left[\left(\begin{array}{c} X_T \\ X_{T,T+1} \cdots X_{T,T+k} \end{array} \right) p_1^{X_{T,T+1}} \cdots p_k^{X_{T,T+k}} \right] \quad (2.12)
 \end{aligned}$$

Questa configurazione presuppone l'indipendenza degli eventi di trasmissione. Ball e Donnelly[11] forniscono una prova che indica che l'indipendenza degli eventi si interrompe quando viene infettata circa la radice quadrata della popolazione sensibile, il che non risulta un problema ai fini del calcolo. È bene ricordare infatti, che la stima

di R_0 deve essere effettuata nel primissimo periodo di un infezione. Riorganizzando i termini la verosimiglianza si riduce ad una Poisson:

$$L(R_0, \mathbf{p}) = \prod_{t=1}^T \frac{e^{-\mu_t} \mu_t^{N_t}}{N_t!} \quad (2.13)$$

dove $\mu_t = R_0 \sum_{j=1}^{\min(k,t)} N_{t-j} p_j$.

A questo punto si può procedere tramite semplici tecniche di massima verosimiglianza per stimare R_0 e p_j .

Dato che in questo caso di studio l'intervallo seriale è già noto non c'è bisogno della costruzione di una massima verosimiglianza che vada a stimare il parametro \mathbf{p} . L'analisi in questo lavoro svolta infatti è stata eseguita dopo che si sono presentati molti focolai nel mondo ed è quindi stato possibile tenere tracce tra i contatti tra casi primari e secondari per la stima dell'intervallo seriale con una certa precisione (errore misurato dagli intervalli di confidenza). In questi caso, l'interesse si concentra quindi sulla sola stima di R_0 .

Di seguito, si descrive lo stimatore per R_0 che risulta molto semplice da implementare, infatti può essere derivato dalla massima verosimiglianza (MLE) 2.13.

Definendo il termine comunemente chiamato score (punteggio) come il gradiente, cioè il vettore delle derivate parziali del logaritmo della funzione di verosimiglianza, ne risulta lo score dalla 2.13 la seguente equazione:

$$U_{R_0}(T) = \sum_{t=1}^T \frac{(N_t - \mu_t)}{R_0} \quad (2.14)$$

Ove $\mu_t = R_0 \sum_{j=1}^{\min(k,t)} N_{t-j} p_j$.

Uguagliandolo a zero e risolvendo per R_0 si ottiene il seguente stimatore di massima verosimiglianza:

$$R_0 = \frac{\sum_{t=1}^T N_t}{\sum_{t=1}^T \sum_{j=1}^{\min(k,t)} p_j N_{t-j}} \quad (2.15)$$

Tramite questo si può procedere al calcolo del valore di R_0 .

2.5.2 R_0 utilizzando campione di osservazioni cinesi e italiane

Senza ripetere quanto detto all'inizio del paragrafo 2.4.4 perchè il procedimento e le considerazioni effettuate sono le medesime si osservano direttamente i risultati, ottenuti questa volta tramite massimizzazione di R_0 nell'equazione 2.13 (in sostanza si utilizza lo stimatore 2.15):

$$R_0^{Cina} : 3.00[2.98, 3.03]^{12}$$

$$R_0^{Italia} : 3.10[3.04, 3.17]^{13}$$

Considerazioni:

- Il metodo usato precedentemente sembra che sovrastimi i parametri
- Anche in questo caso ma come ci si poteva aspettare, l' R_0 italiano appare leggermente superiore a quello cinese e come spiegato in precedenza è plausibile che sia "vero". Essendo poi davvero vicini i due valori e essendo state prese misure diverse (meno restrittive) il risultato è del tutto possibile
- Anche in questo caso si ricorda che si è utilizzato come generation-time, una log-normale (equazione 2.1), con parametri fissi e quindi facendoli variare secondo gli IC la distribuzione cambia e di conseguenza R_0 varia

Un'ultima analisi è volta a capire se i numeri trovati sono significativi. Per farlo basta guardare le stime di R_0 svolte da OMS e altri ricercatori. Facendolo risulta che le stime ufficiali hanno un intervallo di confidenza con questi valori[12]:

$$R_0 = [2, 4] \quad (2.16)$$

L'incertezza su questi valori rimane assoluta come ben spiegato in precedenza. Il risultato trovato ci indica R_0 medio per tale malattia e trasmette un'idea sulla forza di diffusione del coronavirus.

Questo risultato ci indica la pericolosità della malattia. La stima trovata non indica alta pericolosità che invece viene espressa dalla sua velocità di diffusione che può portare a sua volta al collasso delle strutture ospedaliere. È meglio pensare a questo numero ottenuto come un punto di partenza per il comportamento del virus in assenza di fattori umani o ambientali del mondo reale. Per interpretare e capire come nel tempo varia questo parametro in base ai decreti emessi, si utilizza R_t . I governi infatti utilizzano stime di R_0 giornaliere per capire come il virus si stia diffondendo. Tale stime sono, per quanto sofisticate, altamente teoriche. È una stima basata su altre stime, e quindi i margini d'errore sono ancora più ampi. Tuttavia, è una delle uniche metriche per capire che tipo di restrizioni sulla popolazione applicare.

¹²periodo di tempo considerato dal 23/01/2020 al 07/02/2020

¹³periodo di tempo considerato dal 22/02/2020 al 08/03/2020:

Capitolo 3

I modelli epidemiologici

3.1 Modello *SIR*

3.1.1 In termini teorici

Il modello SIR è stato utilizzato per la prima volta da Kermack e McKendrick ed è stato successivamente applicato a una varietà di malattie, in particolare a malattie dell'infanzia con immunità permanente al recupero, come morbillo, parotite, rosolia e pertosse. S , I e R rappresentano il numero di individui sensibili, infetti e recuperati e $N = S + I + R$ è la popolazione totale.

Il modello SIR viene utilizzato laddove gli individui si infettano direttamente (quindi non attraverso un vettore di malattia come ad esempio un animale) e il virus non manifesta periodo di incubazione. Un individuo che guarisce viene considerato immune, ovvero non può più essere infettabile. Il contatto tra le persone è modellato per essere del tutto casuale.

Il tasso di infezione delle persone è proporzionale al numero di persone infette e al numero di persone suscettibili. Se ci sono molte persone infette, le probabilità che una persona sensibile venga a contatto con qualcuno che è infetto sono alte. Allo stesso modo, se ci sono pochissime persone suscettibili, le probabilità che una persona sensibile entri in contatto con un infetto sono inferiori (poiché la maggior parte del contatto sarebbe tra le persone non sensibili-infette).

Il modello SIR può essere riassunto tramite il seguente grafico.

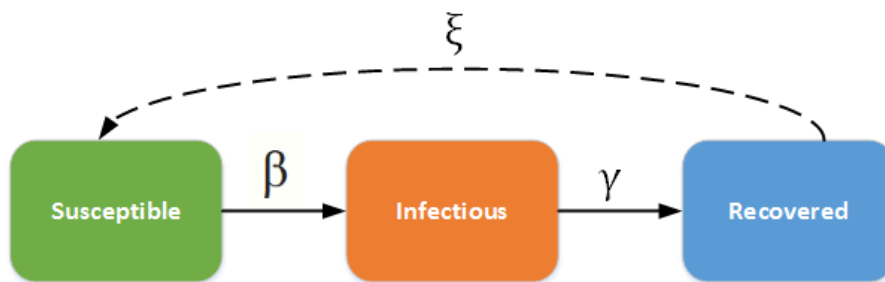


Figura 3.1: Modello SIR

Dove β rappresenta il tasso infettivo e controlla il tasso di diffusione. Rappresenta la probabilità di trasmissione della malattia da un soggetto infetto ad un individuo appartenente alla categoria dei suscettibili. Il tasso di recupero è determinato dalla durata media, D , dell'infezione. Per il modello SIRS, una tipologia differente del modello SIR, ξ è la velocità con cui gli individui guariti, ritornano alla stato sensibile a causa della perdita dell'immunità. Nel caso del coronavirus l'immunità è quasi permanente (si spiegherà meglio nel prossimo modello) e quindi questo aspetto non verrà affrontato.

Nella sua forma deterministica, il modello *SIR* può essere scritto tramite le seguenti equazioni differenziali:

- **Suscettibili:**

$$\frac{dS}{dt} = -\frac{\beta SI}{N} \quad (3.1)$$

- **Infetti:**

$$\frac{dI}{dt} = \frac{\beta SI}{N} - \gamma I \quad (3.2)$$

- **Guariti:**

$$\frac{dR}{dt} = \gamma I \quad (3.3)$$

Questo sistema non è lineare, tuttavia è possibile derivarne la soluzione analitica in forma implicita¹.

¹alcuni strumenti numerici includono i metodi Monte Carlo

Si osserva che:

$$\frac{dS}{dt} + \frac{dI}{dt} + \frac{dR}{dt} = 0$$

da cui ne segue che

$$S(t) + I(t) + R(t) = N \quad (\forall t)$$

Il modello, come facilmente intuibile, è molto limitato. Le ipotesi alla sua base sono molte come spiegato in precedenza, tra cui l'assunzione che la popolazione in esame deve essere in uno stato chiuso. Questo significa che, oltre a non dover esserci nessun entrata/uscita di individui dalla popolazione di riferimento, non ci devono essere né morti né nascite durante il periodo di studio. In una popolazione chiusa senza dinamiche vitali, un'epidemia alla fine si estinguerà a causa di un numero insufficiente di soggetti sensibili per sostenere l'andamento delle infezioni. Ogni persona appartenente al compartimento dei suscettibili verrà sicuramente infettata in una popolazione chiusa. Le persone infette appartenenti a intervalli di generazioni successivi non inizieranno un'altra epidemia a causa dell'immunità permanente della popolazione esistente.

Un esempio semplificativo per la spiegazione di tale fenomeno può essere trovata immaginando un contenitore (paese di riferimento) di particelle² (individui). Essendo lo spazio limitato e le particelle in continuo movimento, è inevitabile che tutte entreranno in contatto fra di loro o in maniera diretta o per conto di terzi.

Le dinamiche di un'epidemia però, come ad esempio l'influenza, sono spesso molto più veloci delle dinamiche di nascita e morte, quindi queste sono spesso omesse in semplici modelli compartimentali, in quanto non significative. Il caso del Covid-19 può essere trattato come una semplice influenza, anche se il tasso di mortalità è leggermente superiore, e quindi questa omissione di dinamicità non ne causa grossi errori.

Il modello SIR è un modello statistico fondamentale in quanto rimodellando le equazioni differenziali si possono determinare altre tipologie di modelli, più complicati, come ad esempio il modello SEIR che tiene conto del periodo di incubazione.

3.1.2 Relazione R_0 con il numero di infetti

Nel capitolo precedente si è trovata una stima di R_0 fondamentale per l'applicazione di questo tipo di modelli. Nel precedente paragrafo nelle equazioni differenziali non essendo presente quest'ultimo non è stato chiarito il motivo della sua importanza all'interno del modello.

²si assuma che abbiano carica neutra per evitare attrazioni e repulsioni

La prima cosa da fare è quella di prendere l'equazione differenziale 3.2, che descrive il numero di infetti nel tempo, e di porla maggiore a zero. In questo modo quello che ne deriva è che un'epidemia avrà luogo se e soltanto se saranno presenti infetti al tempo t :

$$\frac{\beta SI}{N} - \gamma I > 0$$

che con banali operazioni algebriche diventa:

$$\frac{\beta SI}{N} > \gamma I$$

Ora $\frac{S}{N} = 1$ in quanto all'inizio dell'epidemia $S \approx N$ e I si semplifica risultando:

$$\frac{\beta}{\gamma} > 1$$

Ponendo poi una semplice condizione la relazione è trovata.

$$R_0 = \frac{\beta}{\gamma} \tag{3.4}$$

Si fa notare che come descritto al capitolo due abbiamo trovato ora, tramite semplici passaggi algebrici, che un'epidemia si svilupperà (ipotesi iniziale ponendo infetti > 0) solo e soltanto se $R_0 > 1$.

3.1.3 Applicazione modello SIR in Lombardia

Sebbene il modello SIR non consideri un tempo di incubazione per la malattia, che in questo caso è presente, è utile vedere una sua applicazione per capire la portata dell'infezione. Ci si renderà così conto di come lo scenario avrebbe potuto modificarsi se il tempo d'incubazione del coronavirus fosse stato relativamente piccolo o inesistente. Tralasciando quindi decessi e nascite si possono definire i parametri iniziali del modello ricordando che \hat{R}_0 in Italia è risultato pari a 3.12 e 3.60. In questo caso essendo il parametro molto probabilmente una sottostima si è scelto di prendere il valore massimo assumibile dall'intervallo di confidenza trovato nel precedente capitolo. I parametri quindi necessari all'interno del modello sono i seguenti:

- $R_0 = 4$
- γ che è data da $\frac{1}{D}$ ove D numero di giorni della durata media dell'infezione dalla comparsa dei sintomi, per il coronavirus pari a 10[13]

- β è facilmente calcolabile tramite relazione 3.4
- *Parametri iniziali*: in questo si è scelto di inizializzare il modello con i seguenti parametri:
 - $S = N - 1$
 - $I = 1$ (rappresenta il paziente zero)
 - $R = 0$

La figura sottostante riporta il risultato ottenuto. L'asse delle x rappresenta il tempo ossia il numero dei giorni da quando l'epidemia è iniziata mentre l'asse delle y rappresenta il numero percentuale di persone per ogni categoria in un dato istante temporale.

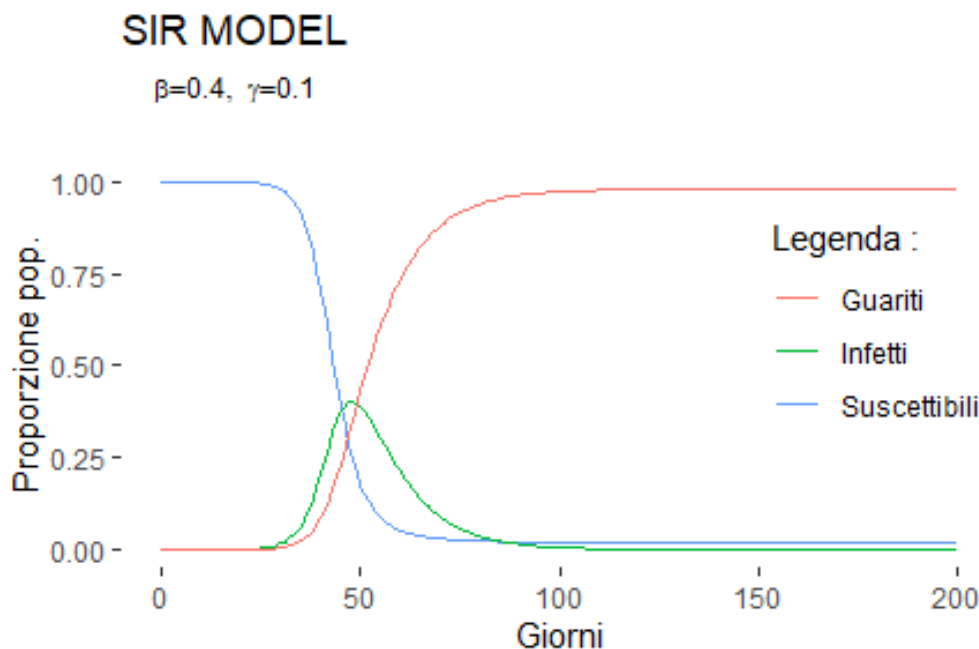


Figura 3.2: SIR Lombardia

L'interpretazione del grafico per ogni linea è la seguente:

- **Linea blu.** Il declino lento della curva, che rappresenta le persone che ancora non sono state infettate, indica la velocità di diffusione della malattia. In questo caso si può concludere che avviene abbastanza lentamente. Intorno al 80 giorno si rileva il valore quasi più basso per tale categoria. Il risultato ci conferma che la lenta diffusione dei contagi all'inizio dell'epidemia, ha reso difficile la comprensione della vera velocità di diffusione del Covid-19

- **Linea verde.** È essenzialmente la curva epidemica data dalla prevalenza (affrontata nel capitolo 1) della malattia. Raggiunge il picco velocemente al giorno 50 circa per poi cadere lentamente nella seconda parte. Il motivo della discesa più lenta è che ci sono più infetti che persone sane e quindi risulta più difficile che una persona contragga il virus. La pendenza della linea verde riflette il tasso di recupero.
- **Linea rossa.** Rappresenta il numero di persone tolte dalla simulazione perché sono guarite e hanno sviluppato l'immunità. La malattia si può definire "risolta" quando tale numero è maggiore dei casi infetti, e quindi quando la linea rossa interseca quella verde. In questo caso circa al 55 giorno.

Quindi se si immagina l'assenza di un periodo di incubazione al Covid-19 e l'assenza di casi asintomatici il risultato che ne fuoriesce dal modello è che il picco di infetti viene raggiunto intorno al 50 giorno.

3.1.4 L'andamento di R_t in Lombardia

In questo paragrafo si cerca di stimare l'andamento del parametro R_0 nel tempo (formalmente R_t), per quanto riguarda la Lombardia, per avere una concezione migliore di come i contagi si siano evoluti nel tempo. La stima sarà derivata considerando gli infetti in funzione del tempo, perché come affermato nella sezione 3.1.2 gli individui malati sono strettamente collegati al basic reproduction number. In particolare gli infetti al tempo t possono essere trovati tramite la seguente relazione:

$$I(t) \approx I_0 e^{(R_0-1) \gamma t} \quad (3.5)$$

Ove $I(t)$ infetti al tempo t e I_0 infetti al tempo iniziale della malattia.

Per stimare R_0 la prima cosa da capire è quale modello utilizzare per aver le migliori performance possibili. In questo caso quelli presi in considerazione sono regressione lineare, regressione log-lineare e regressione poissoniana. La stima del modello, per osservare le variazioni del parametro, è stata effettuata su una finestra mobile di 5 giorni che è sufficiente per stimarne un trend che risulterà crescente o decrescente.

Da un punto di vista grafico, mettendo in relazione gli infetti e il tempo, la miglior soluzione sembra essere un modello di tipo log-lineare, come si può osservare nell'immagine presente nella pagina successiva, in cui si sono raggruppate osservazioni per una finestra temporale di 10 giorni³.

³non si è scelto 5 giorni per non rendere incomprensibile il grafico, il risultato è pressoché equivalente

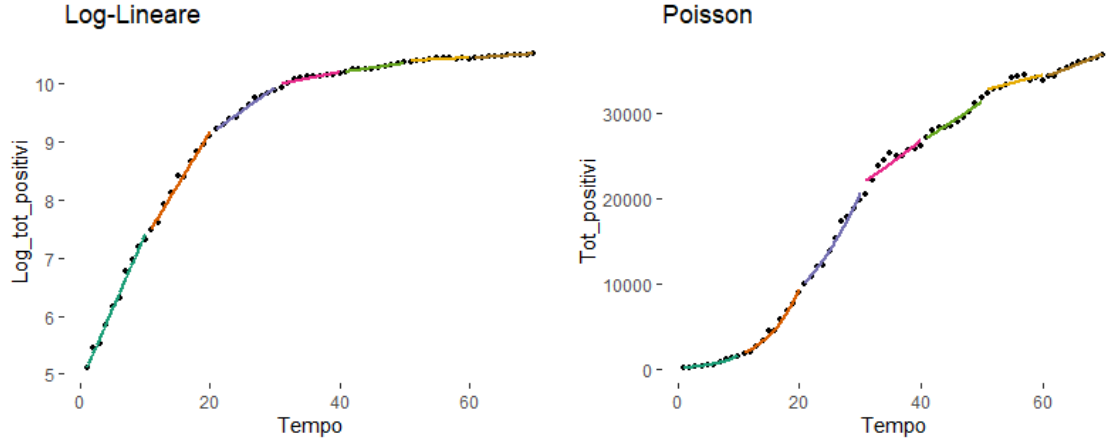


Figura 3.3: Differenze adattamento modelli ai dati

A questo punto siccome il parametro R_0 è strettamente collegato al numero degli infetti e in questo caso non risulta di nessun interesse la costruzione di un modello per la previsione (lo scopo è stimare R_t) si è scelto come metodo per il confronto tra Poisson, regressione lineare e log-lineare l'errore quadratico medio (comunemente abbreviato con MSE). Più tale valore è basso maggiore è l'adattamento del modello ai dati. I tre candidati sono quindi i seguenti modelli:

- Regressione lineare:

$$I(t) = \beta_0 + \beta_1 t + \epsilon_t \quad (3.6)$$

- Regressione di Poisson:

$$\log(\mu_t) = \underset{n \times k}{X'_t} \underset{k \times 1}{\beta} \quad (3.7)$$

ove X'_t matrice trasposta contenente n osservazioni per k modalità e β sono i k coefficienti di regressione.

- Regressione log-lineare:

$$\log(I(t)) = \beta_0 + \beta_1 t + \epsilon_t \quad (3.8)$$

Dato che ogni modello è stato utilizzato su finestre temporali di 5 giorni la scelta è stata effettuata sulla media dei vari valori MSE. I risultati sono i seguenti:

Modello	MSE
Regressione Lineare	233996.2
Regressione Log-Lineare	0.001192147
Regressione di Poisson	0.001202085

Il modello log-lineare risulta, anche se già intravisto graficamente, il migliore e si procede quindi con questo modello.

Applicando il logaritmo nella 3.5 si ottiene questa equazione:

$$\log I(t) \approx \log I_0 + (R_0 - 1) \gamma t \quad (3.9)$$

Preferendo una regressione log-lineare, il modello che si viene a creare per la stima degli infetti è quello definito dall'equazione 3.8. Ove chiaramente i β rappresentano i coefficienti di regressione, β_0 l'intercetta, β_1 è il coefficiente angolare che indica la pendenza della retta e il parametro t è la variabile indipendente o regressore. ϵ_t rappresenta l'errore ed è definito da una variabile casuale che in questo caso assume la forma di una normale.

Se si confrontano le due relazioni, 3.8 e 3.9 si può notare che:

$$\begin{aligned} \hat{\beta}_0 &= \log I_0 \\ \hat{\beta}_1 &= (\hat{R}_0 - 1) \gamma \end{aligned} \quad (3.10)$$

Sistemando la 3.10 tramite alcuni elementari passaggi algebrici si ottiene la seguente relazione:

$$\hat{R}_0 = 1 + \frac{\hat{\beta}_1}{\gamma} \quad (3.11)$$

Ora, data la relazione 3.4 si ottiene un'ulteriore uguaglianza:

$$\hat{R}_0 = 1 + \frac{\hat{\beta}_1}{\gamma} = \frac{\hat{\beta}}{\gamma} \quad (3.12)$$

da cui si può esplicitare beta per ottenere:

$$\hat{\beta} = \left(1 + \frac{\hat{\beta}_1}{\gamma}\right) \gamma = \gamma + \hat{\beta}_1 \quad (3.13)$$

Questa relazione è fondamentale. Si è riusciti a trovare un legame tra il modello lineare, per la stima degli infetti, e il parametro β . Risulta quindi ovvio che trovando tale valore si riesca a stimare anche il parametro R_t .

Si è fatta anche un'ulteriore considerazione. Siccome il numero degli infetti è altamente correlato con il numero di tamponi effettuati giornalmente si è deciso di inserire come peso tale informazione all'interno del modello di regressione log-lineare. In particolar modo ogni unità statistica avrà peso pari a $\frac{1}{N_{\text{tamponi}}}$ in modo da dar peso maggiore se i tamponi effettuati sono pochi. Basti pensare che più tamponi vengono effettuati, più è facile trovare individui positivi e quindi è corretto fare tale osservazione.

Il risultato della stima di R_t è espresso dal grafico 3.4.

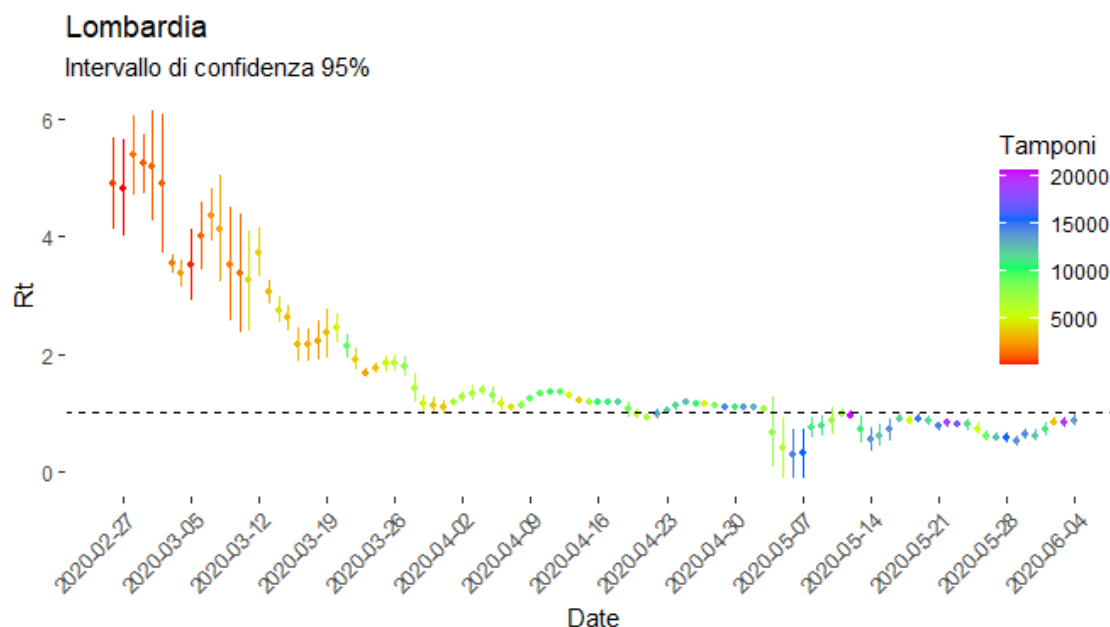


Figura 3.4: Andamento effettivo R_t in Lombardia

Appare evidente un trend decrescente indice di azione governative, come i decreti, funzionanti. Si nota anche che il valore scende sotto la soglia di uno, espressa dalla linea tratteggiata di colore nero, in procinto della fase 2. Le linee prolungate indicano gli intervalli di confidenza, ovvero i valori assunti dal parametro con una certezza in questo caso del 95%. I primi giorni l'andamento è molto più variabile rispetto agli ultimi, questo è indice di un maggior imprecisione nella stima di R_t .

3.2 Modello *SEIR*

3.2.1 In termini teorici

Il modello SEIR è un derivato del modello SIR. S, I e R rappresentano ancora una volta il numero degli individui sensibili, infetti e guariti. Molte malattie hanno una fase latente durante la quale l'individuo è infetto ma non ancora infettivo. Questo ritardo tra l'acquisizione dell'infezione e lo stato infettivo può essere incorporato nel modello SIR aggiungendo una popolazione latente/esposta, E, e lasciando che gli individui infetti (ma non ancora infettivi) passino da S a E e da E a I. Tale lasso di tempo in una malattia altro non è che il periodo di incubazione.

Secondo i dati attualmente disponibili per quanto riguarda il Covid-19, le persone sintomatiche sono la causa più frequente di diffusione del virus. È ritenuto anche possibile che persone nelle fasi prodromiche⁴ della malattia, e quindi con sintomi assenti o molto lievi, possano trasmettere il virus. Tale motivo che molto probabilmente rappresenta la causa della diffusione della malattia, non potrà non essere preso in considerazione.

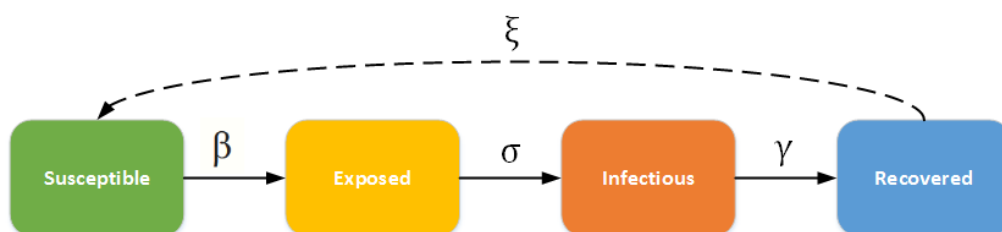


Figura 3.5: Modello SEIR

Il diagramma SEIR /SEIRS, immagine 3.5, mostra come gli individui si muovono attraverso ogni compartimento del modello. La linea tratteggiata mostra come il modello SEIR diventa un modello SEIRS (suscettibile - esposto - infettivo - recuperato - suscettibile), in cui le persone guarite possono diventare nuovamente sensibili (il recupero non conferisce immunità per tutta la vita). Il tasso di velocità con cui accade è dato come nel modello SIR dal valore di ϵ . La malaria ad esempio è una malattia con lunghe durate di incubazione e in cui il recupero conferisce solo un'immunità temporanea. Nel caso del coronavirus non si sa ancora se dopo la malattia permane un'immunità protettiva indotta dagli anticorpi neutralizzanti e quale dose di anticorpi è necessaria perché ciò avvenga. Tutti i dati fino ad ora raccolti, e il

⁴manifestazione morbosa, senza carattere specifico, che precede l'insorgenza dei sintomi caratteristici di una malattia

confronto con un virus parente stretto, il Sars1⁵ suggeriscono che almeno per qualche mese, se non anni, chi guarisce definitivamente non rischia di infettarsi di nuovo[14]. Per questo motivo tale parametro lo si assume pari a 0. Il significato di β e γ rimane identico a quello del caso SIR. Il nuovo parametro σ è il tasso di infezione degli individui latenti ed è dato dal reciproco dei giorni di incubazione del virus.

Nella sua forma deterministica il modello SEIR può essere scritto tramite le seguenti equazioni differenziali:

- **Suscettibili:**

$$\frac{dS}{dt} = -\frac{\beta SI}{N} \quad (3.14)$$

- **Esposti:**

$$\frac{dI}{dt} = \frac{\beta SI}{N} - \sigma E \quad (3.15)$$

- **Infetti:**

$$\frac{dI}{dt} = \sigma E - \gamma I \quad (3.16)$$

- **Guariti:**

$$\frac{dR}{dt} = \gamma I \quad (3.17)$$

Anche in questo caso come nel modello SIR si presuppone che la popolazione sia chiusa, e come già spiegato questa approssimazione non porta a errori grossolani. In questo modello poiché la latenza ritarda l'inizio del periodo infettivo dell'individuo, la diffusione secondaria da un individuo infetto avverrà in un secondo momento rispetto a un modello SIR, che non ha ritardi. Pertanto, un incubazione più lunga comporterà una crescita iniziale più lenta dell'epidemia. Il modello inoltre non includendo la mortalità ha come numero riproduttivo di base, $R_0 = \frac{\beta}{\gamma}$, che non cambia rispetto al modello visto in precedenza.

⁵responsabile dell'epidemia da coronavirus del 2003, simile al virus Mers, epidemia del 2014 nei Paesi arabi

3.2.2 Applicazione del modello in Lombardia

Applicando il modello SEIR nel territorio Lombardo si prenderà atto di come le cose sarebbero potute andare (qualsiasi modello non è indice di verità assoluta) se non fossero stati emessi dei decreti atti a limitare i contagi. Tralasciando quindi morti e nati possiamo definire i parametri iniziali del modello che sono:

- $R_0 = 4$ come in precedenza, per gli stessi motivi spiegati prima
- γ sarà data ancora una volta da $\frac{1}{D}$ ove D numero di giorni della durata media dell'infezione dalla comparsa dei sintomi, per il coronavirus pari a 10.
- β è facilmente ancora calcolabile tramite relazione 3.4
- σ è dato da $\frac{1}{G}$ ove G giorni di incubazione del virus. Tale periodo per chi viene contagiato dal virus SARS-CoV-2 è mediamente di 5 giorni[15]. Dato che il modello implica che nel periodo di incubazione non si è infettivi, a tale numero bisogna sottrarre due giorni. Studi hanno dimostrato che mediamente un individuo è infettivo negli ultimi due giorni prima della comparsa dei sintomi. Si avrà quindi un D pari a $10 + 2 = 12$ e un $G = 3$
- *Parametri iniziali:* in questo caso si è scelto di iniziare il modello con i seguenti parametri uguali a quelli precedentemente usati:
 - $S = N - 1$. A tal proposito si ricorda che la popolazione è chiusa e quindi si presuppone che il primo infetto sia appartenente alla regione Lombardia (non si sa esattamente quanto sia veritiera questa informazione)
 - $I = 1$ rappresenta il paziente zero
 - $R = 0$

Il risultato che ne deriva è visibile dal seguente grafico:

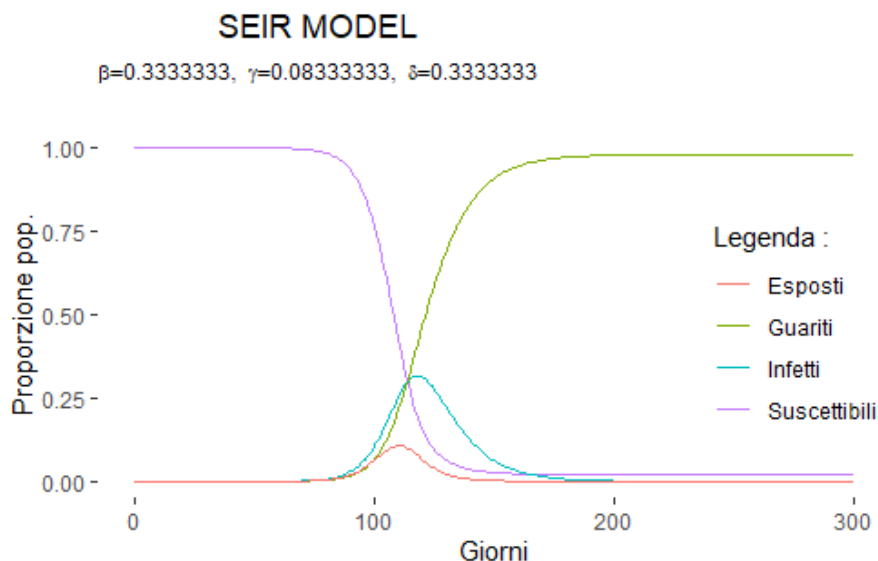


Figura 3.6: SEIR Lombardia

Essendo il modello utilizzato una rappresentazione rispetto al SIR più reale si possono fare interpretazioni più significative:

- **Linea viola.** Indica la velocità di diffusione della malattia. Si nota che appena iniziano a comparire un numero ragionevole di infetti la velocità di diffusione diventa molto rapida
- **Linea blu.** Si raggiunge il picco di infetti al giorno 110 circa per poi cadere lentamente nella seconda parte. Il motivo della discesa più lenta è il medesimo del modello SIR
- **Linea arancione.** Rappresenta la curva degli esposti
- **Linea rossa.** Rappresenta il numero di persone tolte dalla simulazione perchè sono guarite e hanno sviluppato immunità

La curva epidemica, rappresentata da questo modello ci indica che presumibilmente il picco si raggiunge al 110° giorno circa, ovvero ben 4 mesi dopo il primo contagiato. Tale risultato se confrontato con il picco dei valori reali ci mostra come in realtà la malattia sia arrivata sul territorio ben prima del 22 febbraio (data dei primi contagi). Sebben non si possa esprimere con certezza tale risultato risulta davvero impensabile che il picco possa essere stato raggiunto in soli due mesi, anche perchè le restrizioni

prese sono servite per appiattare la curva e quindi a prolungare il più possibile l'arrivo del punto di massimo assoluto. Se anche i parametri del modello vengono modificati, ad esempio abbassando l' R_0 ad una soglia più vicina a 3 (presumendo che non ci sia sottostima dai calcoli del capitolo 2), o aumentandolo in maniera significativa non si riesce a spiegare un picco in due mesi. E' davvero veritiero pensare che la malattia sia in circolo in Lombardia già dalla metà di dicembre. Alcuni medici dei più importanti ospedali Lombardi affermano di aver avuto polmoniti "strane" attorno appunto al 15/17 dicembre, il che porterebbe validità ai risultati espressi dal modello SEIR.

I risultati sono seppur molto interessanti, ancora deboli, in quanto se si tiene ad esempio in considerazione l'ipotesi di popolazione chiusa nella realtà non è del tutto realistica. Una popolazione chiusa studiata dal modello SEIR prevede che gli individui siano sempre messi in contatto gli uni con gli altri portando ad un contagio totale. Questa ipotesi non è del tutto sbagliata, ma chiaramente non aiuta a capire il vero possibile andamento dei contagi. Se ci si pensa, la popolazione in Lombardia è distribuita in maniera del tutto casuale e non si trova in un "contenitore" in cui tutti gli individui vengono a contatto sempre. In maniera più ampia, ad esempio un'analisi sulla diffusione mondiale, un modello SEIR sarebbe del tutto inadeguato. Non è detto che tutti gli individui debbano per forza riscontrare la malattia, come avviene in questo modello!

Un'ultima parte debole è che in questo caso non tutti gli individui che si infettano sono sintomatici, ma ci sono anche individui nella nostra popolazione che risultano asintomatici. L'osservazione non è banale in quanto all'interno del compartimento sono considerati come possibili infettori ma non rientrano negli infetti (I).

L'ultima osservazione interessa l'immunità di gregge che è un fenomeno per il quale, quando un'elevata quota della popolazione risulta immune dall'infezione, gli individui suscettibili sono indirettamente protetti. La proporzione della popolazione immune necessaria per fermare la trasmissione della malattia è in funzione di R_0 (non di R_t) e può essere ricavata dalla seguente formula:

$$1 - \frac{1}{R_0} = 1 - \frac{1}{4} = 75\% \quad (3.18)$$

Il passo finale del lavoro consiste quindi in una costruzione di un modello statistico molto più realistico che cerca di spiegare l'evoluzione della malattia in base agli spostamenti delle persone tramite mezzi pubblici all'interno della regione Lombardia. Non risolveranno chiaramente tutti i problemi descritti fino ad ora ma tale visualizzazione permetterà di trovare risultati più chiari.

Capitolo 4

La struttura dei mezzi pubblici in Lombardia

4.1 La matrice OD

Lo scopo di questo paragrafo è quello di introdurre i dati OD (origine-destinazione) che sono una componente importante per l'applicazione di molti modelli di pianificazione dei trasporti. Ci si servirà di tali informazioni per raggiungere l'obiettivo ultimo che consiste nella costruzione di un modello stocastico che possa mettere in relazione gli spostamenti tramite mezzi pubblici con la diffusione dell'epidemia. Si svolgerà anche un'analisi generale che sarà assestante rispetto al lavoro fatto fino ad ora.

Come suggerisce il nome, i dati origine-destinazione (OD) rappresentano i movimenti attraverso lo spazio geografico, da una certa origine (O) a una destinazione (D). Il nome di matrice le viene attribuito per il fatto che ad ogni paese di origine (j) presente sulle righe viene attribuita una colonna (i) con una destinazione, e l'elemento in posizione ji rappresenta il numero di spostamenti rilevati. Di solito però i dati sono, per facilità di utilizzo, in forma di semplice dataset. All'interno contengono quindi dettagli negli spostamenti tra due punti geografici o, più comunemente, zone¹. La maggior parte dei dati OD fanno riferimenti a posizioni di inizio e fine con colonne, come in questo caso, identificate da degli "*ID*", contenenti stringhe di caratteri che identificano i paesi di origine e destinazione. Alcune volte le stringhe vengono sostituite da latitudine e longitudine delle città di riferimento.

¹spesso rappresentate da un centroide, alcune volte solo dal nome della città/paese

I dati OD in genere contengono più attributi non geografici. È sempre incluso il numero di viaggi che hanno luogo da una certa origine ad una destinazione in un determinato periodo di tempo (ad esempio una tipica giornata lavorativa). Ulteriori attributi possono includere la suddivisione in base alle modalità di trasporto utilizzate per i viaggi. Di solito, come in questo caso, viene catturata una sola modalità. Ad esempio se i viaggi vengono effettuati da una combinazione di modalità bici-treno-camminata viene conteggiata solamente la modalità di viaggio dominante. Ulteriori disaggregazioni dei conteggi complessivi possono comprendere conteggi di viaggio in periodi di tempo diversi (esempio suddivisione per orari).

4.1.1 La matrice OD Lombarda

La matrice origin-destination[16] è un dataset composto da 9 milioni di records e 45 variabili/colonne. Le informazioni presenti altro non sono che un aggiornamento della matrice OD riferita al 2014/2015. Tale dataset è articolato in tabelle per la descrizione sintetica degli spostamenti, con elementi disposti su più righe e su più colonne. Origini e destinazioni sono aggregate in zone di mobilità (sia interne al territorio della Lombardia, sia esterne). La matrice OD 2014 include:

- 8 modalità² (auto conducente, auto passeggero, TPL gomma³, TPL ferro⁴, moto, bici, piedi e altro⁵)
- 5 motivi di spostamento (lavoro, studio, occasionale, affari e rientri a casa)

Per la costruzione di tale dataset nel documento associato all'elaborazione dei dati sono presenti alcune considerazioni:

- La matrice OD-2014 è il risultato della complessa interazione tra modellazioni trasportistiche, questionari on-line, interviste, analisi di indagini disponibili e della domanda esistente rilevata
- I dati si riferiscono ad un giorno feriale medio con periodo da febbraio a maggio
- Gli spostamenti a piedi, sono relativi ai tratti superiori a 10 minuti

²intese come modo prevalente dello spostamento. Qualora uno spostamento fosse costituito da più segmenti, la matrice OD-2014 lo interpreta come unico segmento, che ha come modalità quella prevalente tra i differenti segmenti

³bus urbani, bus extraurbani, pullman e filobus

⁴treno, metro e tram extraurbano

⁵compreso aereo

- Gli spostamenti descritti considerano la modalità prevalente (spostamenti più significativi)
- La zona di riferimento è il singolo comune
- La fascia di popolazione considerata è di età maggiore o uguale a 14 anni

Al fine di rappresentare al meglio la mobilità del territorio lombardo, l'area di studio è stata suddivisa in 1450 zone di mobilità, ove 1264 sono costituite da singoli comuni, 108 da aggregazione di alcuni comuni e 78 costituite dalla disaggregazione di grossi comuni. Nel dataset quindi sono presenti per ogni origine e destinazione dei valori trovati tramite complessi procedimenti che indicano il numero medio di persone che si muovono per causa e mezzo in un giorno lavorativo nel periodo indicato in precedenza.

La definizione della matrice regionale 2014 è basata su un modello di domanda di trasporto, ovvero uno strumento che, alimentato con dati rappresentativi della popolazione e del territorio regionale e calibrato con informazioni rilevate direttamente sul campo, ha lo scopo di rappresentare nel migliore dei modi la realtà. Per poter procedere alla determinazione delle matrici origine-destinazione, è stato utilizzato un approccio complesso a stadi successivi, precisamente 4, mediante opportuni modelli di:

- **Emissione ed attrazione degli spostamenti.** Sono stati utilizzati indici di mobilità. Per emissione si intendono gli spostamenti emessi da un certo paese e per attrazione gli spostamenti in entrata in una certa località. Il modello prevede poi un bilanciamento iterativo del numero di spostamenti attratti con quello degli spostamenti emessi (o viceversa, in funzione del motivo specifico). L'output di questa fase è rappresentato dal totale degli spostamenti emessi ed attratti da ciascuna zona e costituisce la base di partenza per il modello di distribuzione
- **Distribuzione degli spostamenti.** È stato utilizzato un modello gravitazionale che ha ripartito gli spostamenti in funzione di due fattori, attrattività delle zone di destinazione e impedenza temporale tra la zona di origine e quella di destinazione. Utilizzando il grafo di rete stradale, per ciascuna coppia origine-destinazione sono state determinate le matrici delle distanze (in km) e dei tempi (in minuti). Questa fase, inoltre, ha consentito di verificare la connessione di tutte le zone tramite i grafi

- **Ripartizione modale.** Per ogni spostamento, effettuato tra le zone origine e destinazione, per ogni motivo, è stata ricavata l'aliquota di spostamenti che utilizza un certo modo di trasporto. Nello specifico, è stato utilizzato un approccio "kite" considerando separatamente
 - Ripartizione modale degli spostamenti aventi origine in un singolo comune/aggregazione di comuni e destinazione in altro comune/aggregazione di comuni attraverso l'utilizzo del modello *logit/regressione logistica*
 - Ripartizione modale all'interno della diagonale della matrice (spostamenti con origine e destinazione all'interno dello stesso comune/aggregazione di comuni) e, nelle sottomatrici, spostamenti con origine e destinazione all'interno dello stesso comune (per i comuni disaggregati in più zone) tramite un approccio *deterministico*
- **Assegnazione alla rete tramite l'uso di grafi stradali e ferroviari.** Per assegnare la domanda alla rete, viene simulata l'interazione tra
 - domanda di trasporto, rappresentata dalla matrice origine-destinazione della modalità privata e da quella pubblica ferroviaria
 - offerta di trasporto, rappresentata dal grafo di rete stradale e dal grafo dei servizi ferroviari. Si determinano così i flussi di utenti e le prestazioni per ciascun elemento del sistema di offerta.

4.2 Analisi della matrice

Un'analisi sul dataset iniziale è utile per capire le abitudini di spostamento dei cittadini Lombardi. Le prime importanti osservazioni che si possono fare sono le seguenti:

- Gli spostamenti che quotidianamente interessano la Lombardia sono pari a 16,40 milioni e sono cresciuti, rispetto al 2002, del 4,6%
- Il Comune di Milano, come si poteva immaginare, si conferma la città con il maggior utilizzo dei mezzi pubblici

Risulta interessante comprendere, quale sia il mezzo di trasporto più utilizzato all'interno della regione. Il grafico a torta rappresenta la percentuale di impiego di ogni categoria di trasporto.

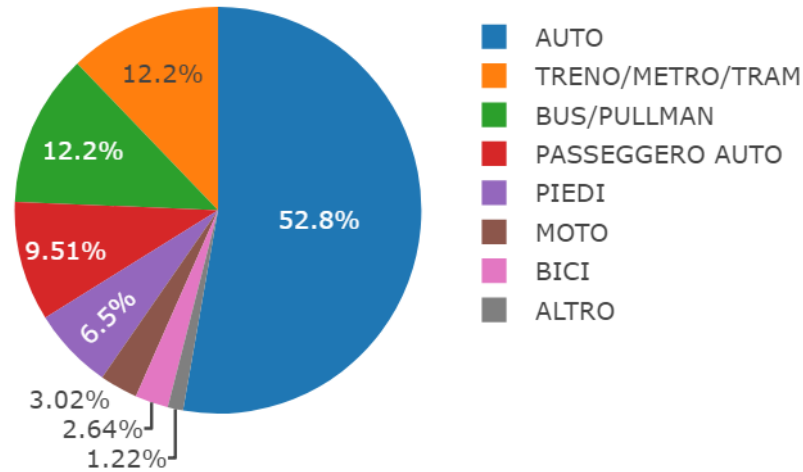


Figura 4.1: Mezzi di trasporto più utilizzati in Lombardia

Si evince come la macchina sia il mezzo di trasporto più utilizzato per più della metà degli spostamenti giornalieri. Se si considerano anche gli individui che si spostano con tale mezzo, senza essere conducenti, si spiegano addirittura il 62% dei movimenti. Treno/metro e tram hanno all'incirca gli stessi individui rispetto ai bus e pullman. Una percentuale minore è rappresentata da spostamenti a piedi in moto e in bici.

Una seconda analisi viene effettuata per capire quali siano i maggiori motivi da cui derivano gli spostamenti giornalieri. Nel grafico sulle ascisse troviamo i motivi di spostamento che in base al colore, dal più scuro (meno intenso) al più chiaro indicano la percentuale rispetto al totale degli spostamenti.

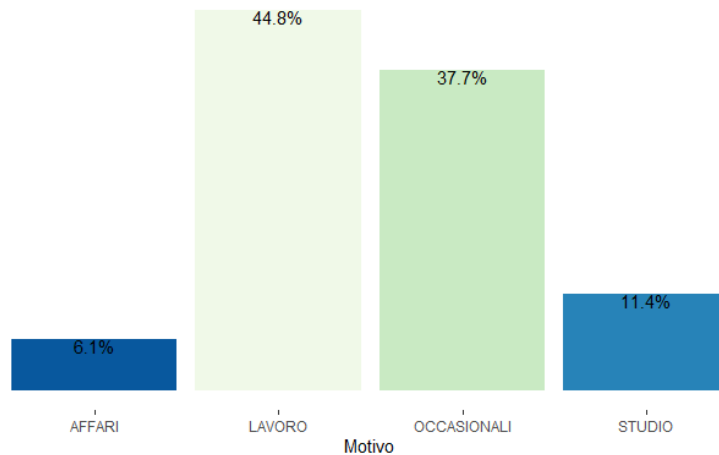


Figura 4.2: Motivi di spostamento in Lombardia

Il lavoro costituisce il motivo principale di spostamento, se unito agli spostamenti con affari si spiega il 50% del motivo di mobilità giornaliera. Una percentuale importante è comunque dovuta agli spostamenti occasionali. Nel grafico le somme delle percentuali non sono pari a 100 in quanto non sono considerati i rientri nelle varie abitazioni.

Un'ultima analisi è fatta per capire all'interno di ogni *motivo di spostamento* quale sia il mezzo più usato. In questo caso nell'asse delle ordinate sono presenti i vari motivi di movimento condizionati al mezzo preso e il colore ne indica l'intensità. Si riportano nel grafico soltanto percentuali, per ovvi motivi, maggiori o uguali a 1.

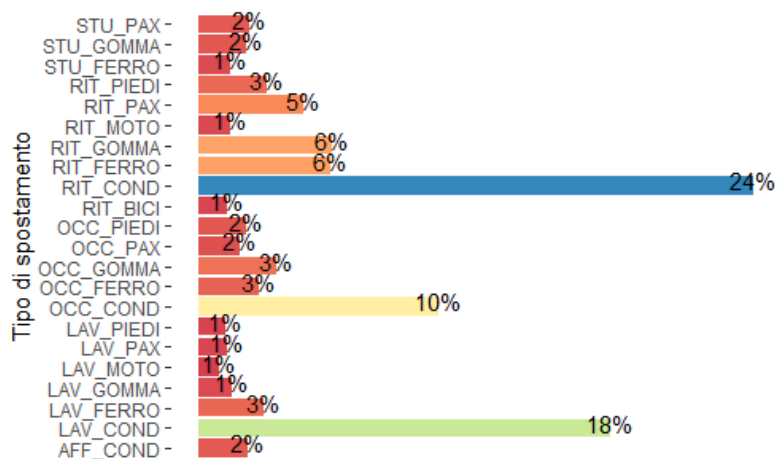


Figura 4.3: Mezzi di trasporto più utilizzati in Lombardia per motivo

Si evince che per ogni motivo di spostamento il mezzo più utilizzato per muoversi è sempre l'auto. I mezzi pubblici costituiscono percentuali al quanto minori, soprattutto per quanto riguarda i motivi di spostamento lavorativi, mentre assumono valori più alti per movimenti occasionali. Emerge anche che gli studenti preferiscono spostarsi maggiormente o come passeggeri in auto o tramite mezzi pubblici.

4.3 I collegamenti dei mezzi pubblici Lombardi

Per quanto riguarda la matrice OD-2019, utilizzata in tale lavoro, costituisce un aggiornamento della OD-2014. Sono stati ad esempio aggiornati i dati relativi a nuove tipologie di collegamenti.

Al fine di costruire la rete dei trasporti dei mezzi pubblici si sono considerate le sole modalità *TPL gomma* e *TPL ferro* indicanti tutti i trasporti pubblici relativi a tutti i possibili motivi di spostamento tranne quello degli affari (spostamenti quasi nulli).

Il risultato è stato il seguente:

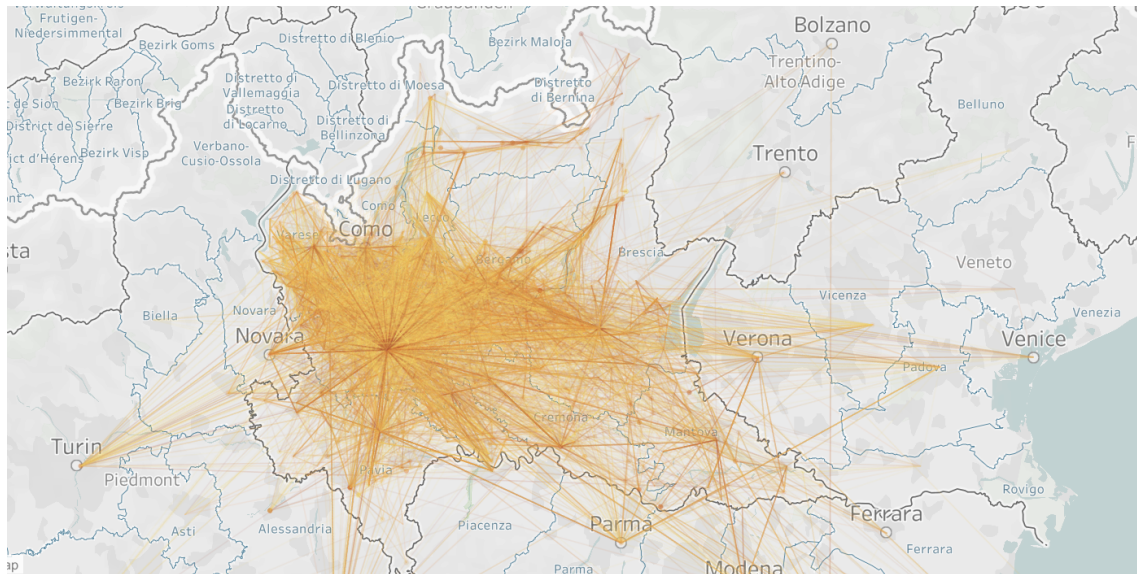


Figura 4.4: La rete dei trasporti dei mezzi pubblici Lombardi

Ove la colorazione più scura indica la presenza di maggiore quantità di persone in quella tratta. È evidente come la rete è al quanto densa e molto concentrata (come ci si poteva aspettare) attorno alla provincia di Milano. Si notano anche collegamenti con altre regioni, verso le quali le persone mediamente ogni giorno si spostano dalla regione Lombardia.

Uno dei tanti problemi nella fase di pulizia è stato quello di gestire comuni come Milano e Monza che nel dataset vengono riportati in micro-aree (ad esempio Milano 1, Milano 2. . .). La decisione è stata quella di aggregare tali informazioni per considerare i vari sottogruppi come unico ente. Le suddivisioni di spostamenti tramite differenti mezzi sono state raggruppate. In questo modo si ha un numero di persone medie che si sposta dal comune A al comune B con qualsiasi mezzo di trasporto. Si ricorda a tal proposito che lo scopo è capire quante persone si spostano da un punto ad un altro e non con che mezzo si muovono.

Capitolo 5

Modello stocastico $SEII_aR$

5.1 Introduzione al modello

Il mondo è in continua urbanizzazione, si stanno formando in questi ultimi decenni ammassi di aree urbane densamente popolate e aree rurali scarsamente popolate. Le persone oggi giorno si spostano continuamente nelle aree urbane e le città stanno crescendo di dimensioni. Trasporti e città affollate sono per i virus come l'ossigeno per il fuoco.

In questo ultimo capitolo si propone una combinazione tra la rete dei trasporti e un modello statistico atto a stimare la diffusione del coronavirus nel territorio Lombardo. Lo scopo è comprendere come l'urbanizzazione e i mezzi di trasporto influenzano la trasmissione delle malattie, nello specifico per il Covid-19.

Alcune considerazioni possono essere già fatte. Nelle aree meno urbanizzate ci saranno meno mezzi di trasporto e quindi quello che ci si aspetta è che in questi luoghi l'andamento delle infezioni sia minore rispetto ad aree più urbanizzate. Si studierà quindi l'interazione tra urbanizzazione e diffusione delle malattie infettive.

I risultati hanno implicazioni dirette per le linee guida e le politiche di controllo delle malattie. Ad esempio è facile immaginare che l'applicazione delle restrizioni ai viaggi ha un impatto maggiore sul numero finale dei malati nei paesi più urbanizzati rispetto ai paesi meno urbanizzati. In particolare come si potrà osservare mediante il modello che si costruirà, le restrizioni ai viaggi sono più efficaci nel ritardare l'epidemia rispetto alla quarantena se applicate subito.

L'urbanizzazione comporta il raggruppamento di persone all'interno di una precisa area geografica. Questa informazione è da tenere in considerazione in quanto più persone si trovano in uno stesso scomparto più in teoria veloce ne sarà la diffusione

della malattia. Detto ciò lo scopo non è quello di descrivere singoli focolai in popolazioni specifiche o di trovare il "miglior modello" o la strategia per un paese specifico, ma piuttosto studiare il fenomeno da un punto di vista più generico per capire come il Covid-19 si sarebbe diffuso se la mobilità non fosse stata limitata.

Si utilizzerà un modello epidemiologico per descrivere l'epidemia di tipo $SEII_aR$ mentre per quanto riguarda la popolazione di ogni singolo comune nel territorio Lombardo, si è utilizzato il dataset di Regione-Lombardia per integrare alla matrice OD-2019 la densità demografica di ogni singolo paese.

Si simulerà la diffusione della malattia infettiva per tutta la regione a partire da un preciso comune.

5.2 La costruzione del modello

5.2.1 Le equazione stocastiche

Il modello di dinamica della malattia è un modello stocastico che può essere descritto come una rete in cui ogni nodo rappresenta una posizione e ogni margine tra le posizioni rappresenta le persone che viaggiano tra i paesi e che quindi possono diffondere ulteriormente la malattia. In ogni luogo esiste una serie separata di equazioni stocastiche che governano la dinamica della malattia locale, al quale si aggiunge un vettore di spostamenti.

Prima di iniziare con termini ed equazioni complesse è bene spiegare il funzionamento del modello da un punto di vista intuitivo tramite la figura 5.1.

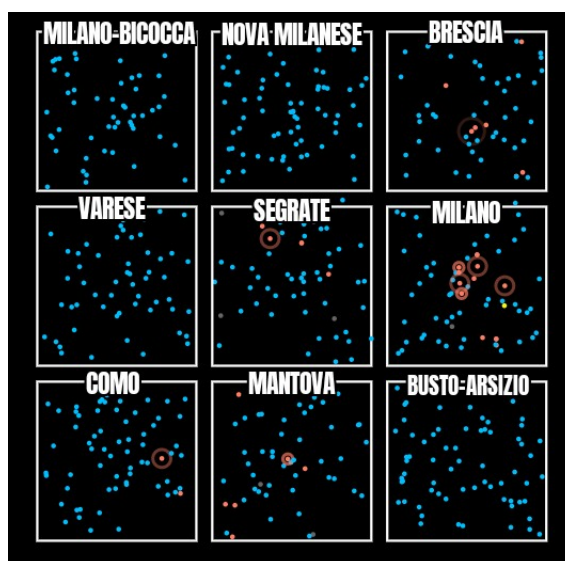


Figura 5.1: Funzionamento modello stocastico di tipo $SEII_aR$

Ogni paese può essere rappresentato tramite dei contenitori che in questo caso hanno nomi scelti in maniera del tutto casuale. Le palline al loro interno rappresentano la popolazione di riferimento. La colorazione rossa dei pallini rappresenta l'individuo infetto mentre rimane azzurra se l'individuo è sensibile. All'interno di ogni contenitore le persone si spostano in modo del tutto casuale, e l'andamento dell'infezione è descritta da un semplice modello, il cui funzionamento rimane identico a quello del capitolo precedente. In questo caso Milano-Bicocca non ha infetti all'interno del suo scomparto e quindi è comprensibile intuire come la malattia non si diffonderà. Avverrà invece il contrario nella città ad esempio di Milano. Eliminando l'ipotesi in cui gli individui sono statici nel loro comune, si permette alle persone di muoversi in diversi scomparti. Si immagini che una persona di Milano che ha contratto il coronavirus (pallina rossa) si sposti per motivi di studio nella città di Milano-Bicocca, allora in tale scomparto comparirà un individuo infetto che potrà contagiare il resto della popolazione.

La dinamica dell'infezione in ogni paese, viene descritta da un modello $SEII_aR$ stocastico. Siano $S^i(t)$ $E^i(t)$ $I^i(t)$ $I_a^i(t)$ il numero di individui sensibili, individui esposti, individui infettivi sintomatici e individui infettivi asintomatici al momento t nella posizione i .

Le equazioni stocastiche di questo modello sono rispettivamente:

- **Suscettibili:**

$$S^i(t + \Delta t) = S^i(t) - \text{Binom}(S^i(t), \frac{\beta \Delta t I^i(t)}{N_i} + P_{ta} \frac{\beta \Delta t I_a^i(t)}{N_i}) \quad (5.1)$$

- **Esposti:**

$$E^i(t + \Delta t) = E^i(t) - \text{Binom}(S^i(t), \frac{\beta \Delta t I^i(t)}{N_i} + P_{ta} \frac{\beta \Delta t I_a^i(t)}{N_i}) - \text{Multinom}(E^i(t), P_a \lambda \Delta t, (1 - P_a) \lambda \Delta t) \quad (5.2)$$

- **Infettiti sintomatici:**

$$I^i(t + \Delta t) = I^i(t) + \text{Binom}(E^i(t), (1 - P_a) \lambda \Delta t) - \text{Binom}(I^i(t), \gamma \delta t) \quad (5.3)$$

- **Infettiti asintomatici:**

$$I_a^i(t + \Delta t) = I_a^i(t) + \text{Binom}(E^i(t), P_a \lambda \Delta t) - \text{Binom}(I_a^i(t), \gamma \delta t) \quad (5.4)$$

Ove β e γ parametri già definiti in precedenza negli altri modelli. N_i rappresenta la numerosità della popolazione nella posizione i , $\text{Binom}(n, p)$ è la distribuzione binomiale con n prove e probabilità di successo pari a p e $\text{Multinom}(n, p_1, p_2)$ è la distribuzione multinomiale con n prove e probabilità di successo pari a p_1 e p_2 . P_{ta} non è altro che la probabilità di trasmissione dell'infettivo asintomatico ridotta rispetto a quello sintomatico mentre P_a è la probabilità di essere un individuo asintomatico. L'equazione per il numero di persone guarite $R^i(t)$ è superflua in quanto in ogni momento $S^i(t) + E^i(t) + I^i(t) + I_a^i(t) + R^i(t) = N_i$ e quindi si può semplicemente ricavarla.

Con questo tipo di modello cambia l'equazione del basic reproduction number che risulta essere:

$$R_0 = \frac{\beta}{\gamma} (P_{ta} P_a + (1 - P_a)) \quad (5.5)$$

Si ricorda che la distribuzione binomiale è una distribuzione di probabilità discreta che descrive il numero di successi in un processo di Bernoulli. Serve a capire il numero di "successi" in n "prove", in questo caso quante persone si infettano ("successo") dati n contatti tra individui. Si fa notare come l'introduzione di questa variabile casuale rappresenti una grossa differenza rispetto ai modelli precedentemente considerati. La distribuzione multinomiale invece è una distribuzione di probabilità discreta che generalizza la distribuzione binomiale in più variabili. In altri termini, laddove la

distribuzione binomiale descrive il numero di successi in un processo di Bernulli, per il quale ogni singola prova può fornire due soli risultati, la distribuzione multinomiale descrive il caso più generale in cui ogni prova possa fornire un numero finito di risultati, ognuno con la propria probabilità[17]. In questo caso serve per attribuire una probabilità di contagio ad un potenziale individuo infetto (che ancora non abbia manifestato sintomi, classe E, esposti), in questo modo:

- p_0 = probabilità di non essere infettato
- p_1 = probabilità di essere infettato e avere sintomi
- p_2 = probabilità di essere infettato senza avere sintomi

5.2.2 Mobilità e dinamica globale delle infezioni

Nel modello, la dinamica della malattia nelle unità di analisi è formulata attraverso gli individui in viaggio. Ogni individuo ha una casa e un luogo di lavoro definiti che identificano rispettivamente l'origine e la destinazione. Durante il giorno, le persone si mescolano nel luogo di lavoro, mentre di notte si mescolano nel luogo di residenza. Si sottolinea l'importanza di tenere traccia degli individui pendolari, vale a dire che bisogna assicurarsi che gli stessi individui siano gli stessi pendolari ogni giorno.

Oltre al pendolarismo, si considerano i viaggi effettuati senza tale dinamica. Il viaggio senza pendolarismo è implementato nel modello consentendo a tutti gli individui che non si spostano tra paesi di viaggiare in una posizione casuale, con una certa probabilità fissa. Il modello distingue tra giorno e notte. Durante il giorno le persone possono infettare/essere infettate nel luogo in cui lavorano, mentre di notte possono infettare/essere infettate nel luogo in cui vivono. Sono gli stessi pendolari che viaggiano avanti e indietro ogni giorno. All'inizio di una giornata, tutti i pendolari vengono inviati al loro posto di lavoro, dove si mescolano per 12 ore. In questo modo i lavoratori possono infettare le persone nella città di destinazione. Una volta finite le 12 ore lavorative i pendolari vengono quindi inviati alle rispettive posizioni di casa, dove si mescolano per le restanti 12 ore. In questo modo nelle 12 ore successive possono infettare le persone della città di origine che a loro volta possono quindi infettare la popolazione di quel paese!

5.3 La scelta dei parametri

La parte fondamentale prima di eseguire il modello è la scelta dei parametri. Le scelte effettuate sono le seguenti:

- P_a che rappresenta la probabilità che un individuo ha di essere asintomatico la si ipotizza pari al 66%[\[18\]](#)
- P_{at} che è la probabilità di trasmissione in meno rispetto ad un indivio sintomatico di un asintomatico la si pone pari al 50% circa [\[19\]](#)
- γ sarà ancora una volta il reciproco della durata media dei giorni dell'infezione che in questo caso è pari a 12 e non 10 per osservazioni già effettuate nel paragrafo 3.2.2
- R_0 lo si pone ancora una volta pari a 4
- β da equazione 5.5 questa volta sarà differente rispetto ai modelli precedenti.

$$\beta = \frac{\gamma R_0}{(P_{ta}P_a + (1 - P_a))} \quad (5.6)$$

- σ cioè il reciproco del periodo di incubazione sarà dato ancora una volta dal reciproco di 3, e non 5
- *Parametri iniziali:* In questo caso la scelta poteva ricadere su qualsiasi comune presente all'interno della regione Lombardia. Si è scelto per differenti ragioni di ipotizzare il primo caso infetto all'interno del comunque di Milano. Uno di questi motivi sono che considerando solo la rete dei trasporti ipotizzare un caso in un comune piccolo farà espandere il contagio in maniera molto lenta il che è molto distante dalla realtà. Per questo si è scelto di considerare quindi questi numeri:

$$\rightarrow S = N_{Milano} - 1$$

$$\rightarrow I_{Milano} = 1 \text{ questo rappresenta il paziente zero}$$

$$\rightarrow I_a = 0$$

$$\rightarrow R_{Milano} = 0$$

5.4 Il modello applicato alla rete Lombarda

Una volta inizializzati i parametri iniziali, si può applicare il modello $SEII_aR$ alla rete dei trasporti Lombardi, immagine 4.1. Il risultato dei contagi è espresso dal seguente grafico che mette in relazione l'avanzare dei giorni con l'incidenza. Per rendere più agevole la visualizzazione dei risultati si è scelto di raggruppare tutti i comuni per la provincia di appartenenza. Si avrà così un'idea generale di come è avvenuta la diffusione. Provando con differenti comuni i risultati cambiano di poco, quello che varia di solito è il tempo con cui si raggiungono i vari picchi. Per quanto riguarda l'asse delle y rappresenta l'incidenza giornaliera mentre per quanto riguarda l'asse delle x rappresenta lo scorrere del tempo, espresso in questo caso in giorni.

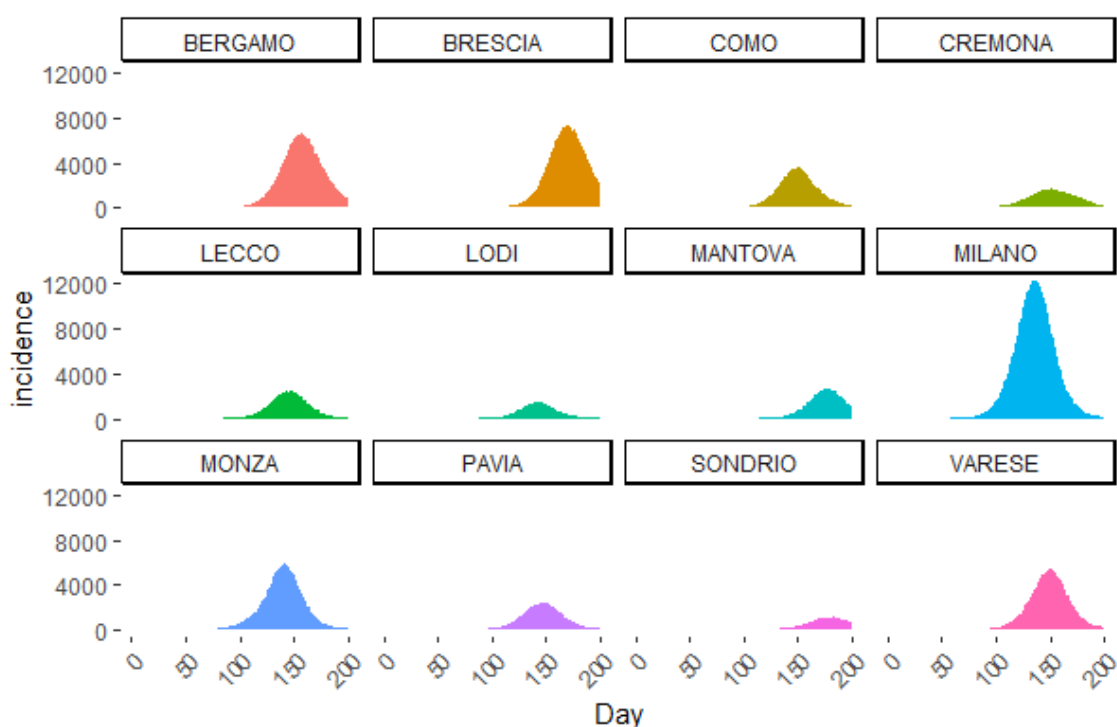


Figura 5.2: Andamento incidenza tramite modello stocastico

Si ha quindi un'immagine di come l'epidemia si sarebbe diffusa se non fossero stati applicati decreti per il contenimento dei movimenti degli individui. La prima evidente osservazione da fare è che la provincia in cui si sarebbe diffusa più velocemente sarebbe stata Milano che è anche la provincia con la più alta densità di popolazione. Il risultato è in linea con le aspettative, i flussi dei trasporti sono maggiormente sviluppati in tale provincia e quindi ogni giorno ci saranno sempre più malati. Si intuisce quindi che la diffusione avviene in base all'intensità degli spostamenti delle

persone. La provincia di Sondrio, che sembra aver una bassa rete di trasporti (si guardi immagine 4.4), è effettivamente quella con la più lenta diffusione. In province invece come Bergamo e Brescia, altamente collegata con Milano, l'andamento dell'infezione è simile anche se arriva con qualche giorno di ritardo. All'inizio dell'epidemia resta evidente come questa si sia sviluppata in maniera lenta, causa forse della sua non individuazione.

Per quanto riguarda la durata del virus si evince ancora una volta come i giorni per l'arrivo del massimo assoluto siano distanti dai 2 mesi avvenuti realmente. Bisogna comunque prendere in considerazione che il picco rilevato nei dati, come già spiegato, non è il picco reale, in quanto non si conosce con esattezza il numero di infetti. Per assurdo si potrebbe aver toccato il massimo anche qualche settimana più tardi senza averne evidenza empirica. Le politiche di intervento del governo sono state fatte in modo tale da appiattire la curva epidemica, ovvero ritardare il picco il più possibile per evitare sovraffollamenti all'interno del sistema sanitario. In questo modo risulta sempre più impossibile che il coronavirus sia presente soltanto da febbraio nel nostro paese. Sarebbe possibile solo se le stime di R_0 fossero completamente sbagliate (servirebbe un valore molto più alto per spiegare due mesi), il che è altamente improbabile. È quindi evidente che si convive con l'infezione almeno dal mese di gennaio, molto probabilmente l'arrivo si spiega intorno alla metà di dicembre ma il modello non esclude neanche il mese di novembre. Purtroppo risulta difficile assegnare un periodo preciso, il massimo che si riesce ad ottenere è quanto detto.

È evidente come limitare l'utilizzo dei trasporti costituisca una grande importanza per la riduzione dei contagi. E' chiaro che tali misure prima vengono applicate maggiore è il loro impatto. Una volta che i contagiati hanno raggiunto e infettato un individuo in un determinato paese, la riduzione dei viaggi non ne costituisce un limite di alcun tipo per la diffusione della malattia in quel paese. Un vincolo degli spostamenti limitati al paese con i primi casi, costituisce il mezzo più efficace per sradicare le infezioni. Nel caso non si riuscisse a procedere in questo modo la quarantena ne risulta il mezzo necessario per abbassare i contagi. Basti pensare di nuovo all'immagine 5.1. Se in quel momento temporale si bloccassero tutti i trasporti in entrata e uscita dalle città di Milano, Brescia, Segrate e Como, non si avrebbero infetti nelle restanti città. Non ci sarebbero infatti vie di accesso per il virus, e quindi quei paesi potrebbe benissimo continuare le attività e gli spostamenti quotidiani.

Conclusioni

In questa tesi sono state trattate molteplici tematiche. Sono stati introdotti diversi tipi di indicatori per lo studio del Covid-19, ognuno con i suoi pregi e difetti. È stata calcolata la curva epidemica di questa specifica malattia per l'Italia e Hubei, fondamentale per la stima del parametro R_0 , necessario poi per la costruzione dei vari modelli, dal più semplice SIR al più complicato $SEII_aR$. Per quest'ultimo sono stati utilizzati i dati relativi agli spostamenti quotidiani delle persone nel territorio Lombardo per capire come si sarebbe evoluto il contagio senza l'introduzione dei decreti.

L'analisi della curva epidemica ha evidenziato come la velocità di diffusione del virus in Italia sia stata maggiore rispetto ad Hubei.

L'analisi dei tassi di mortalità all'interno del territorio italiano ha portato in evidenza che il Covid-19 ha una letalità maggiore di circa 3.5 volte rispetto all'influenza stagionale. I valori trovati sono stati precisamente di 0,056% per il Covid-19 e 0,016% per l'influenza. Questo numero sebbene non risulta eccessivamente elevato, ha indotto a misure restrittive non tanto per la sua pericolosità, bensì per evitare un sovraffollamento delle strutture ospedaliere. Analizzando inoltre le tipologie di persone che sono decedute con la presenza del coronavirus, risulta che il 96% dei casi avevano almeno una malattia pregressa.

La stima del parametro R_0 tramite tasso di crescita esponenziale, il metodo più utilizzato, ha dato come risultato in Italia $\hat{R}_0 = 3.60$ con intervallo di confidenza $[3.53, 3.67]$. La stima tramite massima verosomiglianza, un altro metodo che viene usato, ha portato invece a valori leggermente diversi. La stima puntuale è risultata $\hat{R}_0 = 3.10$ mentre l'IC $[3.04, 3.17]$. I valori con entrambi i metodi sono risultati entrambe le volte più alti rispetto alla stima effettuata sui dati di Hubei. Il motivo potrebbe essere trovato dalle azioni tardive e meno severe all'interno del nostro territorio. Ad oggi non esiste comunque un unico stimatore (i due usati in questo lavoro sono quelli maggiormente utilizzati, ma ne esistono altri) e le molteplici assunzioni per l'utilizzo dei vari metodi quasi mai vengono rispettate nella realtà portando ad imprecisioni. Il risultato ottenuto, che indica quante persone un individuo infetta mediamente, appare più basso di altre malattie come ad esempio Vaiolo $\hat{R}_0 = [3.5, 6.0]$

e Parotite $\hat{R}_0 = [10,12]$. Il problema associato al Covid-19 deriva infatti, non tanto dal valore di R_0 ma dallo studio dell'intervallo seriale che risultando molto più piccolo della Sars ne rappresenta una pericolosità maggiore (sebbene i valori di R_0 sono pressoché identici). Si è scoperto infatti che le infezioni si trasmettono molto velocemente, un individuo può infettarne un altro, a cui possono comparire i sintomi prima della loro manifestazione nella persona trasmittente.

Tramite il valore di R_0 si è trovato inoltre che lo sviluppo dell'immunità di gregge richiede il 75% della popolazione infetta.

Un passaggio ulteriore è stato quello di riuscire, integrando regressione lineare e modello SIR, a trovare l'andamento del *basic reproduction number* giornaliero dopo i decreti, R_t , e a valutarne il suo effettivo trend decrescente prima dell'ingresso nella fase 2. Si è confermato quanto detto dalla protezione civile, ovvero della sua vicinanza a valori come 0.6/0.7 in prossimità del 4 maggio. L'indice sembra anche risalire negli ultimi giorni e se si vuole capire cosa succederà in futuro, bisogna parlare di previsione e intervalli di previsione, un lavoro che non è stato affrontato in questa tesi.

I risultati dei modelli SIR e SEIR hanno prodotto le medesime conclusioni. All'inizio l'epidemia si è innescata in maniera molto lenta (contagi oscurati dal periodo di incubazione e alta presenza di infetti asintomatici) il che può spiegare le manovre tardive effettuate in Italia. La velocità di diffusione cresce in maniera rapida appena buona parte della popolazione diventa infetta. Tutti i modelli inoltre si allontanano dall'ipotesi di insorgenza della malattia nel periodo di febbraio, anticipandola almeno di 1 mese. Si presume dai risultati ottenuti che la malattia sia arrivata nel territorio Lombardo tra la metà di dicembre e la metà di gennaio ma non si esclude un suo approdo nel mese di novembre, essendo il valore di \hat{R}_0 molto variabile.

Dall'analisi del modello stocastico $SEIR_aR$ ne deriva che gli spostamenti tramite mezzi pubblici delle persone facilitano la diffusione del virus. I più colpiti sono i paesi con alta concentrazione di trasporti e popolazione, non a caso l'incidenza più alta è ricaduta sulla provincia di Milano. Si è scoperto anche che meno collegamenti ci sono tra un paese ed un altro più la diffusione avviene lentamente. La densità di popolazione gioca un ruolo chiave nell'andamento degli infetti, in quanto ne è direttamente proporzionale. La riduzione e la chiusura delle reti di trasporto sono un passo fondamentale per fermare e ridurre i contagi se tali misure vengono prese immediatamente. Il motivo è che per far partire un'epidemia in un determinato paese basta una sola persona infetta e lasciando la circolazione libera si ha maggior probabilità che gli individui infetti entrino in molteplici città. Questo modello, più complesso dei due precedenti ha confermato ancora una volta l'impossibilità del raggiungimento del picco d'infetti in un periodo di soli 2 mesi confermando i risultati di SIR e SEIR.

Appendice

LIBRARY

```
1 #Theme set####
2 library(ggplot2)
3 theme_set(theme_classic()+
4             theme(axis.line = element_line(colour =
5               "white"))))
6 #Library####
7 library(dvmisc)
8 library(plotrix)
9 library(incidence)
10 library(ggpubr)
11 library(gridExtra)
12 library(deSolve)
13 library(tidyverse)
14 library(gridExtra)
15 library(spread)
16 library(fhidata)
17 library(viridis)
18 library(gganimate)
19 library(tidyverse)
20 library(plotrix)
21 library(ggribes)
22 library(tidyverse)
23 library(ggplot2)
24 library(ggthemes)
25 library(dygraphs)
26 library(xts)
27 library(tm)
28 library(deSolve)
29 library(wesanderson)
```

CAPITOLO 1

```

1
2 #WORLD SITUATION 24/03/2020#####
3
4 confermati=read.csv("https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid
5 _19_time_series/time_series_covid19_confirmed_global.csv")
6
7
8 ggplot(data = confermati,aes(x = Long, y =
   Lat,fill=X3.24.20)) +
9   borders("legacy_world", colour =
   "grey60",fill="antiquewhite",alpha=0.5) +
10  scale_fill_distiller(type = 'seq', palette =
   "Spectral",direction=1)+
11  theme_map()+
12  geom_point(colour = "black", alpha =
   0.40,shape=21,aes(size=X3.24.20)) +
13  labs(title = "24/03/20
   Infected",fill="Infetti",size="Ampiezza")+
14  theme(panel.background = element_rect(fill = "aliceblue"))
15
16 #ITALY DYNAMIC SITUATION REGION#####
17
18 dat_csv<-read.csv("https://raw.githubusercontent.com
19 /pcm-dpc/COVID-19/master/dati-regioni/dpc-covid19-ita
20 -regioni.csv",header=T)
21 colnames(dat_csv)
22 dat_csv$data=as.Date(dat_csv$data)
23 italy_map <- map_data("italy")
24 print(unique(italy_map$region))
25
26 set.seed(1492)
27 choro_dat <- data_frame(region=unique(italy_map$ region),
28                           value=sample(100, length(region)))
29
30 italy_proj <- "+proj=aea +lat_1=38.15040684902542
31 +lat_2=44.925490198742295 +lon_0=12.7880859375"
32
33 ggplot(data = dat_csv,aes(x = long, y =
   lat,fill=totale_casi,size=totale_casi))+
34  geom_map(data=italy_map, map=italy_map,
35           aes(long, lat, map_id=region),
36           fill="antiquewhite",, size=0.4, color="grey")+

```

```

37   scale_fill_distiller(type = 'seq', palette =
    "Spectral",direction=1)+
38   geom_point(colour = "black", alpha = 0.90,shape=21) +
39   labs(title = "Date: {frame_time}", size =
    "Infetti",subtitle = "INFETTI ITALIA LIVELLO
    REGIONALE") +
40   theme(panel.background = element_rect(fill = "aliceblue"))+
41   transition_time(data)
42
43 #italy dynamic situation province
44
45 dat_csv=read.csv("https://raw.githubusercontent.com/pcm-dpc
46 /COVID-19/master/dati-province/dpc-covid19-ita-province.csv")
47 dat_csv=dat_csv %>%
48   filter(lat!=0 & long !=0)
49 myDateTimeStr1 <- as.vector(dat_csv$data)
50 myDateTimeStr1= gsub("T"," ",myDateTimeStr1)
51 myPOSIXct1 <- as.POSIXct(myDateTimeStr1, format="%Y-%m-%d")
52 dat_csv$data=myPOSIXct1
53 myPoSIXct1 =removeWords(as.character(dat_csv$data),"CET")
54 dat_csv$data=as.Date(myPoSIXct1)
55
56 ggplot(data = dat_csv,aes(x = long, y =
    lat,fill=totale_casi,size=totale_casi))+
57   geom_map(data=italy_map, map=italy_map,
58     aes(long, lat, map_id=region),
59     fill="antiquewhite",, size=0.4, color="grey")+
60   scale_fill_distiller(type = 'seq', palette =
    "Spectral",direction=1)+
61   geom_point(colour = "black", alpha = 0.90,shape=21) +
62   labs(title = "Date: {frame_time}", size =
    "Infetti",subtitle = "INFETTI ITALIA LIVELLO
    PROVINCIALE") +
63   theme(panel.background = element_rect(fill = "aliceblue"))+
64   transition_time(data)
65
66 #EPIDEMIC CURVE ITALY VS HUBEI #####
67
68 cases=read.csv("https://raw.githubusercontent.com/CSSEGISand
69 Data/COVID-19/master/csse_covid_19_data/csse_covid_19_time
70 _series/time_series_covid19_confirmed_global.csv")
71 death=read.csv("https://raw.githubusercontent.com/CSSEGISand
72 Data/COVID-19/master/csse_covid_19_data/csse_covid_19_time
73 _series/time_series_covid19_deaths_global.csv")
74 recovered=read.csv("https://raw.githubusercontent.com/CSSEGIS

```

```

75 andData/COVID-19/master/csse_covid_19_data/csse_covid_19_time
76 _series/time_series_covid19_recovered_global.csv")
77
78 region_confirmed=cases %>%
79   gather(X1.22.20:X4.29.20,key="date",value="n") %>%
80   filter(Province.State=="Hubei" | Country.Region=="Italy" )
81   %>%
82   select(Province.State,date,n)
83
84 region_death=death %>%
85   gather(X1.22.20:X4.29.20,key="date",value="n_morti") %>%
86   filter(Province.State=="Hubei" | Country.Region=="Italy" )
87   %>%
88   select(Province.State,date,n_morti)
89
90 region_recovered=recovered %>%
91   gather(X1.22.20:X4.29.20,key="date",value="n_guariti") %>%
92   filter(Province.State=="Hubei" | Country.Region=="Italy" )
93   %>%
94   select(Province.State,date,n_guariti)
95
96 region=inner_join(region_confirmed,region_death)
97 region=inner_join(region,region_recovered) %>%
98   mutate(casi=n-n_morti-n_guariti) %>%
99   select(Province.State,date,casi)
100
101 region$date=str_replace(region$date, "X", "0")
102 region$date=as.Date(region$date,format = c("%m.%d.%y"))
103 val=which(region$Province.State==c(""))
104 region$Province.State=as.character(region$Province.State)
105 region$Province.State[val]=c("Italia")
106
107 ggplot(region,aes(x=date,y=casi,group=Province.State,
108 fill=Province.State))+
109   geom_density_line(stat = "identity", size=.7,
110     alpha=0.75,colour="Blue4") +
111   scale_fill_manual(values=c("Green4","Yellow3"))+
112   labs(title="Curva epidemiologica",subtitle = "Casi nel
113     tempo",y="Casi",x="Date",fill="Stati")+
114   scale_x_date(breaks="7 days")+
115   theme(axis.text.x = element_text(angle=45,hjust=1))+
116   guides(colour = guide_legend(override.aes = list(size=1)))+
117   theme(legend.justification=c(1,0),
118     legend.position=c(0.6,0.65),
119     legend.text=element_text(size=10))

```



```

114
115 #CINA-->CREATION HUBEI INCIDENCE#####
116
117 data=read.csv("https://raw.githubusercontent.com/CSSEGISand
118 Data/COVID-19/master/csse_covid_19_data/csse_covid_19_time_
119 series/time_series_covid19_confirmed_global.csv")
120 colnames(data)
121
122 infected_a=data %>%
123   filter(Province.State=="Hubei") %>%
124   select(-c("Lat","Long","Country.Region","Province.State"))
125 infected_a=as.numeric(infected_a[1,])
126
127 #Incidence
128 incidence_1=vector()
129 for(i in 2:length(infected_a)){
130   incidence_1[i-1]=infected_a[i]-infected_a[i-1]
131 }
132 incidence_cina=c(incidence_1)
133 incidence_cina=incidence_cina[1:50]
134 time=length(incidence_cina)
135
136 #Incidence output
137
138 date=seq.Date(from = as.Date("2020-01-23"),length.out = time
139   , by = "day")
140 length(date)
141 inc=data.frame(date=date,incidence=incidence_cina)
142 ince<-xts(inc,order.by = date)
143 ince=data.frame(date=ince$date,incidence=incidence_cina)
144 ince$incidence=as.numeric(ince$incidence)
145 ince$date=as.Date(ince$date)
146
147 picchi=ince\%>%
148   summarize(massimo=max(incidence),
149   date=date[which((incidence)==max(incidence))])
150
151 g1=ggplot(ince,aes(x=date,y=incidence))+
152   geom_bar(stat="identity",fill="springgreen3",
153   col="white")+geom_line()+
154   geom_point(picchi,mapping=aes(x=date,y=massimo),
155   col="Red",size=3)+
156   theme(axis.text.x = element_text(angle=45,hjust=1))+
157   theme(legend.position = "None")+
158   labs(title = "Incidenza Cina-Hubei",y="Incidenza")+

```

```

158   scale_x_date(breaks="7 days")+
159   theme(axis.text.x = element_text(angle=45,hjust=1))
160
161 #Incidence with lag 3 days
162
163 input = 1:length(ince$incidence)
164 multiple_of_3 = (input %% 3) == 0
165 valori=input[multiple_of_3]
166
167 somma=function (x,len){
168   aggregate(x,by=list(rep(1:(length(x)/len),each=len)
169   ),FUN=sum)
170 }
171
172 mobile=data.frame(incidence=somma(incidence_cina
173 [-c(49,50)],3),date=ince$date[valori])
174
175 picchi=mobile\%>\%
176   summarize(massimo=max(incidence.x),
177   date=date[which((incidence.x)==max(incidence.x))])
178
179 g2=ggplot(mobile,aes(x=date,y=incidence.x))+
180   geom_bar(stat="identity",fill="springgreen3",col="white")+
181   geom_line()+
182   geom_point(picchi,mapping=aes(x=date,y=massimo),
183   col="Red",size=3)+
184   theme(axis.text.x = element_text(angle=45,hjust=1))+
185   theme(legend.position = "None")+
186   labs(title = "Incidenza Cina-Hubei",subtitle = "Finestra
187   mobile 3 giorni",y="Incidenza")+
188   scale_x_date(breaks="7 days")+
189   theme(axis.text.x = element_text(angle=45,hjust=1))
190
191 ggarrange(g1, g2,ncol=2, nrow=1)
192
193 #ITALY INCIDENCE#####
194
195 data=read.csv("https://raw.githubusercontent.com/CSSEGISand
196 Data/COVID-19/master/csse_covid_19_data/csse_covid_19_time_
197 series/time_series_covid19_confirmed_global.csv")
198 colnames(data)
199 infected=data \%>\%
200   filter(Country.Region=="Italy") \%>\%
201   select(-c("Lat","Long","Country.Region","Province.State"))
202 infected=as.numeric(infected[1,-c(1:29)])

```

```

201 infected
202
203 incidence_italia=vector()
204 for(i in 2:length(infected)){
205   incidence_italia[i-1]=infected[i]-infected[i-1]}
206 length(incidence_italia)
207 incidence_italia
208
209 date=seq.Date(from = as.Date("2020-02-21"), length.out =
    length(incidence_italia), by = "day")
210 length(date)
211
212 inc=data.frame(date=date,incidence=incidence_italia)
213 ince<-xts(inc,order.by = date)
214
215 ince=data.frame(date=ince$date,incidence=incidence_italia)
216 ince$incidence=as.numeric(ince$incidence)
217 ince$date=as.Date(ince$date)
218
219 tamponi=vector()
220 for(i in 2:length(dat_csv$tamponi)){
221   tamponi[i-1]=dat_csv$tamponi[i]-dat_csv$tamponi[i-1]}
222 tamponi
223
224 ince=cbind(ince[-c(1:4)],,tamponi=tamponi)
225
226 picchi=ince\%>\%
227   summarize(massimo=max(incidence),
228   date=date[which((incidence)==max(incidence))])
229
230 g1=ggplot(ince,aes(x=date,y=incidence))+
231   geom_bar(mapping=aes(fill=tamponi),stat="identity",
232   col="white")+geom_line()+
233   theme(axis.text.x = element_text(angle=45,hjust=1))+
234   labs(title = "Incidenza
    Italia",fill="Tamponi",y="Incidenza")+
235   scale_x_date(breaks="7 days")+
236   geom_point(picchi,mapping=aes(x=date,y=massimo),
237   col="Red",size=3)+
238   theme(axis.text.x = element_text(angle=45,hjust=1))+
239   scale_fill_viridis(discrete = F,option = "D")+
240   theme(legend.justification=c(1,0),legend.position=c(0.95,0.4),
241   legend.text=element_text(size=10))
242
243

```

```

244 #Incidence with lag 3 days
245
246 input = 1:length(ince$incidence)
247 multiple_of_3 = (input %% 3) == 0
248 valori=input[multiple_of_3]
249
250 mobile=data.frame(incidence=somma(incidence_italia
251 [-c(76,77,78,79,80)],3),date=ince$date[valori],
252 tamponi=somma(tamponi[-c(76,77,78,79,80)],3))
253
254 picchi=mobile[>]
255     summarize(massimo=max(incidence.x),
256     date=date[which((incidence.x)==max(incidence.x))])
257
258 g2=ggplot(mobile,aes(x=date,y=incidence.x))+
259     geom_bar(mapping=aes(fill=tamponi.x),stat="identity",
260     col="white")+geom_line()+
261     geom_point(picchi,mapping=aes(x=date,y=massimo),
262     col="Red",size=3)+
263     theme(axis.text.x = element_text(angle=45,hjust=1))+
264     labs(title = "Incidenza Italia",subtitle = "Finestra
265     mobile 3 giorni",fill="Tamponi",y="Incidenza")+
266     scale_x_date(breaks="7 days")+
267     theme(axis.text.x = element_text(angle=45,hjust=1))+
268     scale_fill_viridis(discrete = F,option = "D")+
269     theme(legend.justification=c(1,0),
270     legend.position=c(0.95,0.4),
271     legend.text=element_text(size=10))
272
273 ggarrange(g1, g2,ncol=2, nrow=1)
274
275 #Death analysis
276
277 #MORTI#####
278
279 dat_csv<-read.csv("https://raw.githubusercontent.com/pcm-dpc
280 /COVID-19/master/dati-regioni/dpc-
281 covid19-ita-regioni.csv",header=T)
282 myDateTimeStr1 <- as.vector(dat_csv$data)
283 myDateTimeStr1= gsub("T"," ",myDateTimeStr1)
284 myPOSIXct1 <- as.POSIXct(myDateTimeStr1,
285     format="%Y-%m-%d")
286 dat_csv$data=myPOSIXct1
287 myPoSIXct1 =removeWords(as.character(dat_csv$data),"CET")
288 dat_csv$data=as.Date(myPoSIXct1)

```

```

286
287 dat_csv\%>%
288   filter(data=="2020-05-29")\%>%
289   summarize(morti=sum(deceduti),
290     guariti=sum(dimessi_guariti))\%>%
291   mutate(tasso=(morti/guariti))
292
293 dat_csv\%>%
294   summarize(positivi=sum(nuovi_positivi))
295
296 dat_csv\%>%
297   filter(data=="2020-03-10")\%>%
298   summarize(morti=sum(deceduti),
299     guariti=sum(dimessi_guariti))\%>%
300   mutate(tasso=(morti/guariti))
301
302
303 dat_csv\%>%
304   filter(data=="2020-05-10")\%>%
305   summarize(morti=sum(deceduti),
306     guariti=sum(dimessi_guariti))\%>%
307   mutate(tasso=(morti/guariti))

```

CAPITOLO 2

```

1 #R0 ESTIMATION-CHINA#####
2
3 library(R0)
4 china=region%>%
5   filter(Province.State=="Hubei")
6
7 GT = generation.time ("lognormal",c(4.7,2.6))
8 order(infected_a,decreasing = T)
9 est.R0.EG(epid = infected_a, GT, begin =1, end =15,reg.met =
10   "poisson")
11 est.R0.ML(infected_a,GT,begin = 1,end=15)
12
13 #R0 ESTIMATION-CHINA#####
14
15 GT = generation.time ("lognormal",c(4.7,2.6))
16 order(incidence_italia,decreasing = T)
17 #same time china's period
18 est.R0.EG(epid = infected , GT, begin =1, end =15,reg.met =
19   "poisson")
20 est.R0.ML(infected,GT,begin =1,end=15)

```

CAPITOLO 3

```

1 #PARAMETERS CHOICE#####
2
3 R0=4
4 Incubazione=3 ##5 medi - 2 medi dove sono contagioso
5 Periodo_infetto=10 ##10 medi
6 gamma=1/Periodo_infetto #recovered rate
7 sigma=1/Incubazione
8 beta=gamma*R0
9 beta
10
11 #SIR MODEL LOMBARDY#####
12
13 sir <- function(time, state, parameters) {
14   with(as.list(c(state, parameters)), {
15     dS <- -beta * S * I
16     dI <- beta * S * I - gamma * I
17     dR <- gamma * I
18     return(list(c(dS, dI, dR)))
19   })
20 }

```

```

21
22 N_Lombardia=10.04*10^6
23
24 #Input
25
26 init<-c(S =(N_Lombardia-20-1)/N_Lombardia,I
      =20/N_Lombardia,R =0/N_Lombardia)
27 parameters <- c(beta = beta, gamma = gamma)
28 times <- seq(0, 200, by = 1)
29 out <- ode(y = init, times = times, func = sir, parms =
      parameters)
30 out <- as.data.frame(out)
31 out$I=out$I
32 out$R=out$R
33 out$time <- NULL
34 head(out, 10)
35
36 out
37 title <- bquote("SIR MODEL")
38 subtit <- bquote(list(beta==.(parameters[1]),
39 ~gamma==.(parameters[2])))
40
41 sir<-ggplot(out,aes(x=0:200))+
42   ggtitle(bquote(atop(. (title),atop(bold(. (subtit))))))+
43   geom_line(aes(y=S,colour="Suscettibili"))+
44   geom_line(aes(y=I,colour="Infetti"))+
45   geom_line(aes(y=R,colour="Guariti"))+
46   ylab(label="Proporzione pop.")+
47   xlab(label="Giorni")+
48   labs(color="Legenda : ")+
49   theme(legend.justification=c(1,0),
      legend.position=c(1,0.3))
50 sir
51
52 #RO VARIATION IN LOMBARDY####
53
54 detach("package:RO", unload=TRUE)
55 library(tidyverse)
56
57 dat_csv<-read.csv("https://raw.githubusercontent.com/pcm-dpc
58 /COVID-19/master/dati-regioni/dpc-covid19-ita-regioni.csv",
59 header=T)
60 colnames(dat_csv)
61 dat_csv=dat_csv %>%
62   filter(denominazione_regione=="Lombardia")

```

```

63
64 days<-dim(dat_csv)[1]
65 colnames(dat_csv)
66 dat_csv$t<-1:days
67 days
68 myDateTimeStr1 <- dat_csv$data
69 myDateTimeStr1= gsub("T"," ",myDateTimeStr1)
70 myPOSIXct1 <- as.POSIXct(myDateTimeStr1, format="%Y-%m-%d
    %H:%M:%S")
71 days_dy<-as.Date(myPOSIXct1)
72
73 #need to update manually
74 dat_csv_dy<-xts(dat_csv[, -c(1:6,18:20)], order.by = days_dy,
    frequency = 7)
75 length(dat_csv_dy$ricoverati_con_sintomi)
76 gruppi=c(rep(1,10),rep(2,10),rep(3,10),rep(4,10),rep(5,10)
77 ,rep(6,10),rep(7,10))
78 dat_csv_dy=cbind(dat_csv_dy,log=log(as.numeric
79 (dat_csv_dy$totale_positivi)),as.numeric(gruppi))
80
81 p1=ggplot(dat_csv_dy)+
82   geom_point(aes(x=t,y=log),size=1)+
83   scale_color_brewer(palette="Dark2")+
84   geom_line(aes(x=t,y=log,group=gruppi,col=as.factor(gruppi)),
85     stat="smooth",method="lm",se = F,alpha=1,size=1)+
86   theme(legend.position = "none")+
87   labs(y="Log_tot_positivi",x="Tempo",title = "Log-Lineare")
88
89 p2=ggplot(dat_csv_dy)+
90   geom_point(aes(x=t,y=totale_positivi),size=1)+
91   scale_color_brewer(palette="Dark2")+
92   geom_smooth(aes(x=t,y=totale_positivi,group=gruppi,
93     col=as.factor(gruppi)),method="glm",method.args =
94     list(family = "poisson"),se = F,size=1)+
95   theme(legend.position = "none")+
96   labs(y="Tot_positivi",x="Tempo",title="Poisson")
97
98 grid.arrange(p1,p2,ncol=2)
99
100 duration<-15
101 beta_vec<-NULL
102 sd_vec<-NULL
103 mse_POIS=NULL
104 mse_LM=NULL
105 mse_LOGLM=NULL

```



```

105
106 tamponi=vector()
107 for(i in 2:length(dat_csv$tamponi)){
108   tamponi[i]=dat_csv$tamponi[i]-dat_csv$tamponi[i-1]}
109 tamponi[1]=dat_csv$tamponi[1]
110 tamponi[3]=492
111 dat_csv$tamponi=tamponi
112
113 for (i in 3:(days-2)){
114   fit <- glm((totale_positivi)~t,weights =
115     1/tamponi,family="poisson",data=dat_csv[(i-2):(i+2),])
116   mse_POIS=c(mse_POIS,get_mse(fit))
117   fit <- glm((totale_positivi)~t,weights =
118     1/tamponi,family="gaussian",data=dat_csv[(i-2):(i+2),])
119   mse_LM=c(mse_LM,get_mse(fit))
120   fit <- glm(log(totale_positivi)~t,weights =
121     1/tamponi,family="gaussian",data=dat_csv[(i-2):(i+2),])
122   mse_LOGLM=c(mse_POIS,get_mse(fit))
123   beta_vec<-c(beta_vec,coef(fit)[2])
124   sd_vec<-c(sd_vec,coef(summary(fit))[2,2])
125 }
126 mean(mse_LM)
127 mean(mse_LOGLM) #the best
128 mean(mse_POIS)
129
130 label<-as.Date(substr(dat_csv$data,1,10))[3:(days-2)]
131 mean <- 1+(beta_vec*duration)
132 lower <- 1+((beta_vec-1.96*sd_vec)*duration)
133 upper <- 1+((beta_vec+1.96*sd_vec)*duration)
134 df <- data.frame(label, mean, lower,
135   upper,tamponi=tamponi[-c(1,2,length(mean),length(mean)-1)])
136 fp <- ggplot(data=df, aes(x=label, y=mean, ymin=lower,
137   ymax=upper,col=tamponi)) +
138   geom_pointrange(size=0.2) +
139   scale_color_gradientn(colours = rainbow(5))+
140   geom_hline(yintercept=1,
141     size=0.5,color="black",linetype="dashed") +
142   xlab("Date") + ylab("Rt")+
143   labs(title = "Lombardia",subtitle = "Intervallo di
144     confidenza 95%",col="Tamponi" )+
145   scale_x_date(breaks="7 days")+
146   theme(axis.text.x = element_text(angle=45,hjust=1))+
147   theme(legend.justification=c(1,0),
148     legend.position=c(1,0.3))
149 fp

```

```

142
143 #MODELLO SEIR LOMBARDY#####
144
145 seir_model = function (current_timepoint, state_values,
146   parameters){
147   S = state_values [1]          # susceptibles
148   E = state_values [2]          # exposed
149   I = state_values [3]          # infectious
150   R = state_values [4]          # recovered
151   with (
152     as.list (parameters),{
153       dS = (-beta * S * I)
154       dE = (beta * S * I) - (delta * E)
155       dI = (delta * E) - (gamma * I)
156       dR = (gamma * I)
157       results = c (dS, dE, dI, dR)
158       list (results)
159     }
160   )
161 }
162
163 parameter_list = c (beta = beta, gamma = gamma, delta=sigma)
164
165 W = N_Lombardia-1
166 Y = 0
167 Z = 0
168 X = 1
169 N = W + X + Y + Z
170
171 initial_values = c (S = W/N, E = X/N, I = Y/N, R = Z/N)
172 timepoints = seq (0, 300, by=1)
173 output = lsoda (initial_values, timepoints, seir_model,
174   parameter_list)
175
176 title <- bquote("SEIR MODEL")
177 subtit <- bquote(list(beta==.(parameter_list[1]),
178   ~gamma==.(parameter_list[2]),~delta==.(parameter_list[3])))
179 output=as.data.frame(output)
180
181 seir<-ggplot(output,aes(x=time))+
182   ggtitle(bquote(atop(. (title),atop(bold(. (subtit))))))+
183   geom_line(aes(y=S,colour="Suscettibili"))+
184   geom_line(aes(y=I,colour="Infetti"))+
185   geom_line(aes(y=R,colour="Guariti"))+
186   geom_line(aes(y=E,colour="Esposti"))+

```

```

185   ylab(label="Proporzione pop.")+
186   xlab(label="Giorni")+
187   labs(color="Legenda :")+
188   theme(legend.justification=c(1,0), legend.position=c(1,0.15))
189
190 seir

```

CAPITOLO 4

```

1 # OD MATRIX-GENERAL ANALYSIS#####
2
3 matrice_od <- read.csv("C:/Users/matte/Desktop/Lavori
   tesi/File_Excel/matrice_od.csv")
4 od=matrice_od
5
6 val=od%>%
7   filter(ZONA_ORIG!=ZONA_DEST) %>%
8   summarise_each(funs(sum),-c(ZONA_ORIG,PROV_ORIG,PROV_DEST,
9     ZONA_DEST,FASCIA_ORARIA))%>%
10  gather(LAV_COND:RIT_ALTRO,key="Tipo",value="N")
11 for(i in 1:length(val$N)){
12   val$N[i]=round(val$N[i],0)
13 }
14
15 val%>%
16   filter(N>100000)%>%
17   mutate(pct = round(prop.table(N),3))%>%
18   ggplot(aes(x=N,y=Tipo,fill=N))+
19   geom_bar(stat="Identity")+
20   scale_fill_distiller(type = 'seq', palette =
     "Spectral",direction=1)+
21   theme(axis.title.x=element_blank(),
22     axis.text.x=element_blank(),
23     axis.ticks.x=element_blank(),legend.position
     ="NULL")+
24   geom_text(aes(label=scales::percent(pct)), angle = 360,
25     position=position_dodge(width=0.2), vjust=0.2)+
26   labs(y="Tipo di spostamento")
27
28 val$Tipo
29 gruppi=c(rep("LAVORO",8),rep("STUDIO",8),rep("OCCASIONALI",8),
30 rep("AFFARI",8),rep("RIENTRO",8))
31 new_val=cbind(val,gruppi)
32

```

```

33 new_val%>%
34   filter(gruppi!="RIENTRO")%>%
35   group_by(gruppi)%>%
36   summarize(N=sum(N))%>%
37   mutate(pct = prop.table(N))%>%
38   ggplot(aes(x=gruppi,y=N,fill=N))+
39   geom_bar(stat="Identity")+
40   scale_fill_distiller(palette =4)+
41   theme(axis.title.y=element_blank(),
42         axis.text.y=element_blank(),
43         axis.ticks.y=element_blank(),legend.position
44         ="NULL")+
45   geom_text(aes(label=scales::percent(pct)),
46             position=position_dodge(width=0.5), vjust=0.9)+
47   labs(title ="Motivo di spostamento",x="Motivo")
48
49 val$Tipo
50 gruppi=gsub("[^;]*_(.*)", "\\1", val$Tipo)
51 new_val_2=cbind(val,gruppi)
52
53 new_val_2=new_val_2%>%
54   group_by(gruppi)%>%
55   summarize(N=sum(N))
56
57 new_val_2$gruppi=as.character(new_val_2$gruppi)
58 new_val_2$gruppi[3]=c("AUTO")
59 new_val_2$gruppi[4]=c("TRENO/METRO/TRAM")
60 new_val_2$gruppi[5]=c("BUS/PULLMAN")
61 new_val_2$gruppi[7]=c("PASSEGGERO AUTO")
62
63 fig <- plot_ly(new_val_2, labels = ~gruppi, values = ~N,
64               type = 'pie')
65 fig <- fig %>% layout(xaxis = list(showgrid = T, zeroline =
66               FALSE, showticklabels = FALSE),
67               yaxis = list(showgrid = T, zeroline =
68               FALSE, showticklabels = FALSE))
69
70 fig
71
72 #CLEAN OD_MATRIX TO BUILDING PUBLIC TRANSPORT MATRIX#####
73
74 matrix_od=od %>%
75   filter(ZONA_ORIG!=ZONA_DEST) %>%
76   select(ZONA_ORIG,PROV_ORIG,ZONA_DEST,PROV_DEST,LAV_FERRO,

```

```

74   LAV_GOMMA , OCC_FERRO , OCC_GOMMA , RIT_GOMMA ,
75   RIT_FERRO , STU_GOMMA , STU_FERRO) %>%
76   mutate(all=LAV_FERRO+LAV_GOMMA+OCC_FERRO+OCC_GOMMA+
77   RIT_GOMMA+RIT_FERRO+STU_GOMMA+STU_FERRO) %>%
78   select(ZONA_ORIG , PROV_ORIG , ZONA_DEST , PROV_DEST , all) %>%
79   filter(all>0)
80
81   #ripetitions eliminated
82
83   new_file=aggregate(matrix_od$all ,
84     by=list(matrix_od$ZONA_ORIG , matrix_od$ZONA_DEST) , FUN=sum)
85
86   new_file=matrix_od%>%
87     group_by(ZONA_ORIG , ZONA_DEST , PROV_ORIG , PROV_DEST) %>%
88     summarize(all=sum(all))%>%
89     rename(from=ZONA_ORIG , to=ZONA_DEST , n=all)
90
91   #eliminated number in strings
92   new_file$from=as.character(new_file$from)
93   new_file$to=as.character(new_file$to)
94
95   split=function(val){
96     x = unlist(strsplit(val , split = '\\s+'))
97     val=paste0(x[Reduce('|' , lapply(c('[A-Za-z]', '67') ,
98       grepl , x))], collapse = ' ')
99     return(val)
100   }
101
102   for(i in 1:length(new_file$from)){
103     new_file$from[i]=split(new_file$from[i])
104     new_file$to[i]=split(new_file$to[i])
105   }
106
107   new_file$from=as.factor(new_file$from)
108   new_file$to=as.factor(new_file$to)
109   costruzione_province=new_file
110
111   new_file$from=as.character(new_file$from)
112   new_file$to=as.character(new_file$to)
113   new_file=new_file[,c(1:2,5)]
114   new_file=new_file %>%
115     filter(from!=to)
116
117   new_file$from=as.factor(new_file$from)
118   new_file$to=as.factor(new_file$to)

```

```
117  
118 new_file=aggregate(new_file$n,  
119                     by=list(new_file$from,new_file$to), FUN=sum)  
119 colnames(new_file)=c("from","to","n")  
120  
121  
122  
123 #write.csv(new_file,file = "OD.csv")
```

CAPITOLO 5

```

1 #DATAFRAME WITH CITY AND POPULATION#####
2
3 com<-read.csv("C:/Users/matte/Desktop/Lavori_tesi/File_excel
4 /11125_comuni Popolazione-residente-totale-per-classe-di-eta-
5 Totali-al-11---Comun.csv", sep=";")
6 popolazione=as_tibble(com)
7 colnames(popolazione)
8
9 popolazione=popolazione %>%
10   filter(Anno==2019 & Livello.territoriale=="Comune") %>%
11   select(Territorio,Totale) %>%
12   transmute(Territorio=toupper(Territorio),Totale=Totale)
13
14 citta=levels(new_file$from)
15 pop=vector()
16 v=0
17 for(i in citta){
18   v=v+1
19   if(i %in% popolazione$Territorio){
20     a=which(popolazione$Territorio==i)
21     a=a[length(a)]
22     pop[v]=popolazione$Totale[a]}
23   else {
24     pop[v]="Na"
25   }
26 }
27 length(citta)
28 length(pop)
29
30 regioni=c("LOMBARDIA","LAZIO","CAMPANIA","SICILIA","VENETO",
31 "EMILIA-ROMAGNA","PIEMONTE","PUGLIA","TOSCANA","CALABRIA",
32 "SARDEGNA","LIGURIA","MARCHE","ABRUZZO","FRIULI VENEZIA
33   GIULIA",
34 "TRENTINO ALTO ADIGE","UMBRIA","BASILICATA","MOLISE","VALLE
35   D'AOSTA","ASTI","BIELLA","BOLOGNA","BOLZANO","CANTU",
36 "CASSANO D'ADDA","CASTIGLIONE D'ADDA","CUNEO","FERRARA",
37 "FORLICESENA","MODENA","MUGGIO ","NOVARA","PADOVA","PARMA",
38 "PEREGO","PIACENZA","RAVENNA","REVERE","RIMINI","TORINO",
39 "TRENTO","TREVISO","VENEZIA","VICENZA")
40 regioni_pop=c(10018806,5898124,5839084,5056641,5056641,
41 4448841,4392526,4063888,3742437,1965128,1653135,1565307,
42 1538055,1322247,1217872,1062860,888908,570365,310449,126883,
43 76164,44616,388367,106951,40007,18911,4665,56124,132009,

```

```

42 117946,184727,23579,104284,936887,194417,1757,102355,159057,
43 2508,148908,886837,117417,885447,261905,112198)
44
45 for(i in 1:length(regioni)){
46   val=which(citta==regioni[i])
47   pop[val]=regioni_pop[i]
48 }
49
50 for(i in 1:length(citta)){
51   val=which(pop=="Na")
52   pop[val]=1500
53 }
54
55 #STARTED VALUES####
56
57 E=rep(0,length(pop))
58 R=rep(0,length(pop))
59 I=rep(0,length(pop))
60 Ia=rep(0,length(pop))
61 citta_completo=data.frame(location_code=citta,
62 S=as.numeric(pop),
63 E=E,I=I,Ia=Ia,R=R)
64
65 #MODEL####
66
67 #STARTED VALUES FROM MILANO
68
69 which(citta_completo$location_code=="MILANO")
70 citta_completo[830,]
71 citta_completo[830,]$I=1
72 citta_completo[830,]$R=0
73 citta_completo[830,]$E=0
74 citta_completo[830,]
75 citta_completo
76 which(new_file$from=="TOSCANA" & new_file$to=="AFRICA")
77 new_file=new_file[-423,]
78
79 set.seed(123)
80 d <- commuter(
81   seiiar=citta_completo,
82   commuters=new_file,
83   r0=4.0,
84   latent_period = 3,
85   infectious_period = 12,
86   asymptomatic_prob=2/3,

```



```

87   asymptomatic_relative_infectiousness=0.10,
88   days_simulation=200,
89   N=1
90 )
91 a=which(d$location_code=="MILANO")
92 d[a]
93
94 #BUILD A PROVINCE/CITY DATAFRAME####
95
96 location_code=levels(citta_completo$location_code)
97 provincia=vector()
98 v=0
99
100 for(i in location_code){
101   v=v+1
102   a=which(costruzione_province$from==i)
103   a=a[1]
104   provincia[v]=as.character(costruzione_province$PROV_ORIG[a])
105 }
106
107 provincia=replace(provincia, provincia=="BG", "BERGAMO")
108 provincia=replace(provincia, provincia=="BS", "BRESCIA")
109 provincia=replace(provincia, provincia=="CO", "COMO")
110 provincia=replace(provincia, provincia=="CR", "CREMONA")
111 provincia=replace(provincia, provincia=="LC", "LECCO")
112 provincia=replace(provincia, provincia=="LO", "LODI")
113 provincia=replace(provincia, provincia=="MN", "MANTOVA")
114 provincia=replace(provincia, provincia=="MI", "MILANO")
115 provincia=replace(provincia, provincia=="MB", "MONZA")
116 provincia=replace(provincia, provincia=="PV", "PAVIA")
117 provincia=replace(provincia, provincia=="SO", "SONDRIO")
118 provincia=replace(provincia, provincia=="VA", "VARESE")
119
120 province=data.frame(location_code=location_code,
121   provincia=provincia)
122
123 ##MODEL, BUILD WITH PROVINCE#####
124
125 d <- merge(d,province,
126   by.x="location_code",by.y="location_code")
127 county <- d[,.(
128   S=sum(S),
129   E=sum(E),
130   I=sum(I),
131   Ia=sum(Ia),

```

```

131 R=sum(R),
132 incidence=sum(incidence),
133 pop=sum(pop)
134 ),
135 keyby=(provincia,location_code,week,day,is_6pm)]
136
137 #select only lombardy's province
138 lombardia=which(county$provincia=="BERGAMO" |
139   county$provincia=="BRESCIA" | county$provincia=="COMO" |
140   county$provincia=="CREMONA" | county$provincia=="LECCO" |
141   county$provincia=="LODI" | county$provincia=="MANTOVA" |
142   county$provincia=="MILANO" | county$provincia=="MONZA" |
143   county$provincia=="CREMONA" | county$provincia=="PAVIA" |
144   county$provincia=="SONDRIO" | county$provincia=="VARESE")
145
146 county_2=county[lombardia,]
147 colnames(county_2)
148
149 #write.csv2(county_2,file = "infetti lombardia.csv")
150
151 #model output#####
152
153 pop_province=county_2%>%
154   group_by(provincia)%>%
155   filter(day==1)%>%
156   summarize(pop=sum(S))
157
158
159 somme=county_2%>%
160   group_by(provincia,day)%>%
161   summarize(S=sum(S),I=sum(I),R=sum(R),E=sum(E),Ia=sum(Ia))
162
163 popprov=left_join(somme,pop_province)
164
165 somme=popprov%>%
166   mutate(Ss=S/pop,Ii=I/pop,Rr=R/pop,Er=E/pop,Iaa=Ia/sum(Ia))
167
168 ggplot(somme,aes(x=day))+
169   geom_line(aes(y=Ss,colour="Suscettibili"))+
170   geom_line(aes(y=Ii,colour="Infetti"))+
171   geom_line(aes(y=Rr,colour="Guariti"))+
172   geom_line(aes(y=Er,colour="Esposti"))+
173   geom_line(aes(y=Iaa,colour="Asintomatici"))+
174   facet_wrap(~provincia)+
175   ylab(label="Proporzione pop.")+

```

```

170   xlab(label="Giorni")+
171   labs(color="Legenda : ")+
172   theme(legend.justification=c(2,0),
         legend.position=c(1,0.3))
173
174
175   colnames(county_2)
176   p <- ggplot(county_2, aes(x=day, y=incidence))+
177     geom_col()+
178     facet_wrap(~provincia)+
179     labs(x="Giorni",y="Curva incidence")+
180     theme(legend.position = "None")+
181     theme(axis.text.x = element_text(angle=45,hjust=1))
182   p
183   #Save plot
184   ggsave(p,device = "png",width = 10,height = 7,filename =
         "Andamento incidence.png")

```

Bibliografia

- [1] Treccani. Le parole del coronavirus. URL: http://www.treccani.it/magazine/parolevalgono/Le_parole_del_Coronavirus/index.html.
- [2] Farr William. Causes of death in england and wales. *Second Annual Report of the Registrar General of Births, Deaths and Marriages in England*, 2:69–98, 1927.
- [3] Sylvia Richardson. Coronavirus statistics: what can we trust and what should we ignore? *The Guardian*, 2020. URL: <https://www.theguardian.com/world/2020/apr/12/coronavirus-statistics-what-can-we-trust-and-what-should-we-ignore>.
- [4] Regione Lombardia Sistema Socio Sanitario ASST Garda. Esecuzione tamponi covid-19 per accertamento guarigione, 2020. URL: <https://www.asst-garda.it/notizie/tamponi-covid-19-di-controllo/>.
- [5] Istituto Superiore di Sanità (ISS). Caratteristiche dei pazienti deceduti positivi all’infezione da sars-cov-2 in italia. *Science*, 2020. URL: <https://www.epicentro.iss.it/coronavirus/sars-cov-2-decessi-italia>.
- [6] William Ogilvy Kermack and Anderson G McKendrick. A contribution to the mathematical theory of epidemics. *Proceedings of the royal society of london. Series A, Containing papers of a mathematical and physical character*, 115(772):700–721, 1927.
- [7] StatGroup-19. Meme ed epidemiologia, 2020. URL: <https://statgroup-19.blogspot.com/>.
- [8] Andrei R. Akhmetzhanova Hiroshi Nishiuraa, Natalie M. Lintona. Serial interval of novel coronavirus (Covid-19) infections. *International Journal of Infectious Diseases*, 93:284–286, 2020.
- [9] Anderson RM. *Infectious diseases of humans: dynamics and control*. Oxford, may r.m. oxford university press edition, 1991.

- [10] Pagano M. White LF. A likelihood-based method for real-time estimation of the serial interval and reproductive number of an epidemic. *Stat med.*, 27(16):2999–3016, 2008.
- [11] Donnelly P. Ball F. Strong approximations for epidemic models-stochastic processes and their applications, 1995.
- [12] Istituto Superiore di Sanità (ISS). Che cos'è r_0 e perché è così importante. *ISS*, 2020. URL: https://www.iss.it/primo-piano/-/asset_publisher/o4oGR9qmvUz9/content/id/5268851.
- [13] Annals of Internal Medicine. The incubation period of coronavirus disease 2019 (covid-19) from publicly reported confirmed cases: Estimation and application, 2020. URL: <https://www.acpjournals.org/doi/10.7326/M20-0504>.
- [14] Bai-Zhong Liu Quan-Xin Long. Antibody responses to sars-cov-2 in patients with covid-19. *Nature Medicine*, 26:845–848, 2020.
- [15] Johns Hopkins University. New study on covid-19. *Bloomberg School of public health*, 2020. URL: <https://www.jhsph.edu/news/news-releases/2020/new-study-on-COVID-19-estimates-5-days-for-incubation-period.html>.
- [16] Regione Lombardia-Open Data. La matrice di dati Origin Destinazione, 2019. URL: <https://www.dati.lombardia.it/Mobilit-e-trasporti/Matrice-OD2016-Passeggeri/tezw-ewgk>.
- [17] L'enciclopedia libera Wikipedia. Distribuzione multinomiale, 2019. URL: https://it.wikipedia.org/wiki/Distribuzione_multinomiale.
- [18] Sergio Romagnani. Professore ordinario di Immunologia clinica dell'Università di Firenze, 2020. URL: https://www.repubblica.it/salute/medicina-e-ricerca/2020/03/16/news/coronavirus_studio_il_50-75_dei_casi_a_vo_sono_asintomatici_e_molto_contagiosi-251474302/.
- [19] Jeffrey Shaman. Substantial undocumented infection facilitates the rapid dissemination of novel coronavirus (sars-cov-2). *Science*, 368:489–493, 2020.