

Analisi del network dei criminali londinesi

Chiara Camerota

August 24, 2018

Introduzione

Il network preso in esame riassume le relazioni tra i membri delle gang londinesi, basandosi sugli atti criminali commessi in gruppo, durante il periodo 2005-2009. I dati vengono forniti in forma anonima dalla polizia, sulla base degli arresti e delle condanne subite dai membri.

La rete presenta 54 nodi, ognuno corrispondente ad un membro delle gang, mentre gli archi possono assumere diversi valori, corrispondenti a diversi livelli di criminalità. Nello specifico: 1 indica che i nodi si frequentano, 2 indica che i due hanno commesso atti criminali assieme, 3 che hanno commesso gravi atti criminali assieme, 4 indica la condizione precedente e che vi è un legame di parentela tra i due nodi.

Il nostro obiettivo è quello riuscire a capire che ruolo gioca la nazionalità nel commettere atti criminali di gruppo, indi per cui rendiamo binario il nostro network, ovvero diamo pari peso a tutti gli archi.

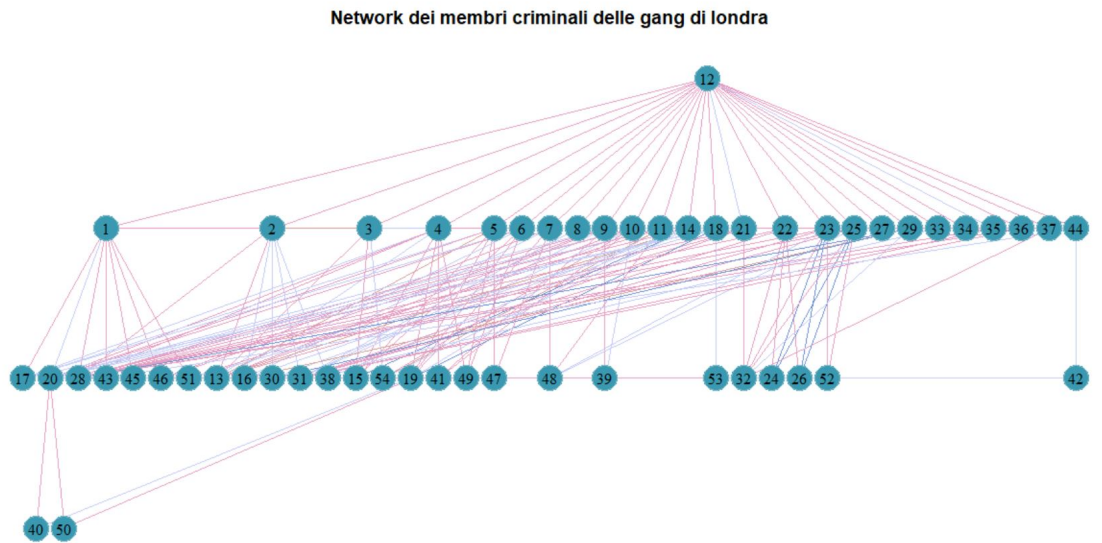


Figure 1: Sopra é riportato il grafo originale, ogni colore indica un legame diverso tra le unità. Il rosa indica un arco di tipo 1, il blu cobalto chiaro un legame di tipo 2, il marrone uno di tipo 3, il blu cobalto scuro un legame di tipo 4 .Come si può osservare, il grafo presenta una struttura piramidale, tipica delle organizzazioni criminali.

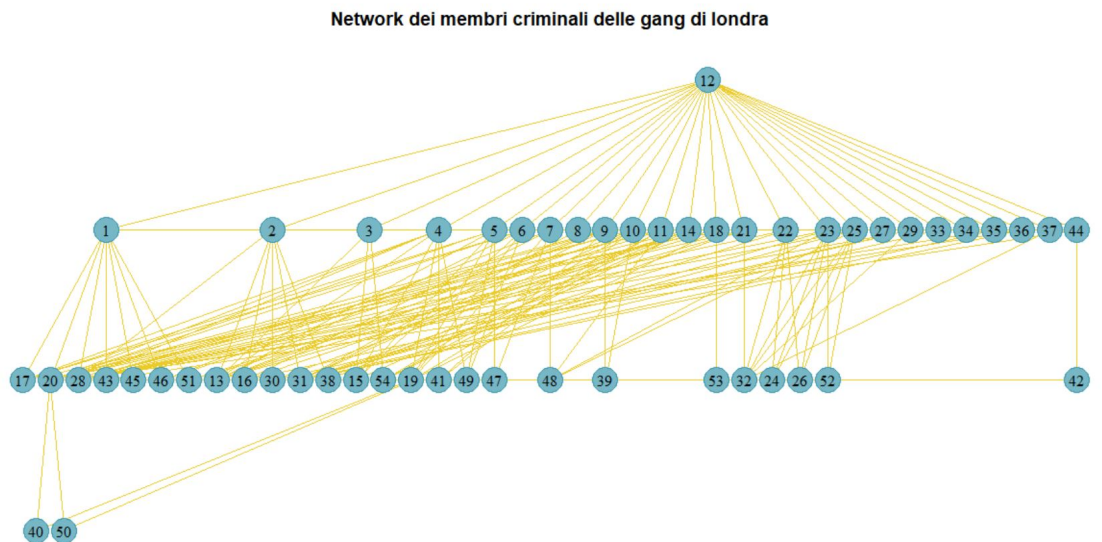


Figure 2: Sopra é riportato il grafo dicotomizzato.

Gli attributi presi in considerazione sono: l'età, il luogo di nascita, la

residenza, il numero degli arresti, il numero di condanne, se il membro é o meno in carcere.

Di seguito sono riassunti in tabella le descrittive delle covariate.

- *Residence*: vale 1 se il criminale ha la residenza, 0 altrimenti

Residence	1	0
Frequenza	20	34

- *Prison*: vale 1 se il criminale é in prigione, 0 altrimenti

Prison	1	0
Frequenza	20	34

- *BirthPlace*: vale 1 per i nati in Africa dell'ovest , 2 per quelli dei Caraibi, 3 per i nati negli UK, 4 per quelli dell'Africa dell'est

BirthPlace	1	2	3	4
Frequenza	12	12	24	6

- *Age*: indica l'età del criminale

Age	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Valore	16.00	18.00	19.00	19.83	21.00	27.00

- *Arrests*: indica il numero di arresti subiti dal criminale

Arrests	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Valore	0.000	5.000	8.000	9.907	14.750	23.000

- *Convictions*: indica il numero di condanne a carico del criminale

Convictions	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Valore	0.000	1.000	3.000	4.204	7.000	13.000

Studi precedenti, come ad esempio Grund, T. and Densley, J. (2012, 2015) dimostrano come sia piú probabile, per due criminali che agiscono insieme, che questi abbino la stessa nazionalità piuttosto che altro in comune; ma vedremo come l'omofilia, ovvero la scelta di commettere atti criminali con persone aventi o meno la propria nazionalità, é molto pi'u forte nelle triadi, piuttosto che nelle 2-star. Una possibile spiegazione di questo fenomeno é che per due criminali (aventi la stessa nazionalità) é piú semplice trovare un amico in comune all'interno di questo gruppo, piuttosto che all'esterno.

Chapter 1

Analisi descrittive e Studio della centralit 

1.0.1 Densit , Censimento delle triadi, transitivit  ed assortative mixing

Iniziamo ora l'analisi descrittiva della rete, come prima cosa possiamo notare che la rete non direzionata ha 54 nodi collegati tra di loro mediante 315 archi, sui $\binom{54}{2} = 1431$ possibili, conseguentemente la probabilit  di osservare un arco (ovvero la densit )   pari a 0.22, essa non   molto alta, ma questo   giustificato dalla natura clusterizzata della rete.

$$\rho = \frac{\text{num. archi osservati}}{\text{num. archi totali}} = 0.22$$

Essendo il network non direzionato non ha senso censire le diadi, ma risulta interessante censire le triadi, in quanto spiegano la struttura gerarchica dei dati.

Tipo di relazione	Frequenza assoluta
nessuna	12528
$A \rightleftharpoons B, C$	9032
$A \rightleftharpoons B \rightleftharpoons C$	2384
$A \rightleftharpoons B \rightleftharpoons C, A \rightleftharpoons C$	860

Come   possibile aspettarsi, la relazione $A \rightleftharpoons B \rightleftharpoons C$   pi  presente della relazione $A \rightleftharpoons B \rightleftharpoons C, A \rightleftharpoons C$, quindi favorisce legami di tipo 2-star, piuttosto che quelli di tipo triangolo. Ovvero   pi  semplice trovare legami che percorrono la piramide verticalmente (per esempio: legami tra capi o legami tra sottoposti), piuttosto che legami triangolari (per esempio: legami tra sottoposti e diversi capi).

Per capire meglio questi tipi di legame, ci viene in supporto la transitivit , ovvero la probabilit  di osservare un legame tra due unit  quando

queste hanno già un nodo in comune.

$$C = \frac{\#triangoli}{\#triangoli + \#2 - star} = 0.52$$

Essendo questo indice poco informativo in forma cruda, é conveniente normalizzarlo.

$$\tau = \log \left(\frac{\frac{C}{1-C}}{\frac{\rho}{1-\rho}} \right) = \log \left(\frac{C}{1-C} \right) - \log \left(\frac{\rho}{1-C\rho} \right) = 1.343943$$

Essendo $\tau > 1$ possiamo confermare, quello che fino ad ora poteva risultare solo una congettura, ovvero che la struttura della rete é gerarchica.

Passiamo ora ad analizzare le covariate in relazione alla rete, ovvero allo studio dell'assortative mixing, di seguito sono riportate in tabella gli assortative coefficient per ognuna di questa che corrisponde alla modularity normalizzata ,per gli attributi qualitativi, e la covariance measures normalizzata, per quelli quantitativi.

Attributo	Assortative coefficient
<i>Birthplace</i>	0.1021293
<i>Residence</i>	0.003307925
<i>Prison</i>	0.0148802
<i>Age</i>	0.1516753
<i>Arrests</i>	0.07194567
<i>Conviction</i>	0.0971864

Possiamo notare come per *BirthPlace* e *Age* sembri esserci, anche se minima, una tendenza positiva a commettere reati insieme a parità di questi attributi (presenza di omofilia); per gli altri tratti caratterizzanti non sembra esserci un effetto di omofilia. In altre parole, i gruppi criminali, sembrano scegliere i componenti anche basandosi su l'età e il luogo di origine.

1.0.2 Analisi della centralità

Concentriamoci ora sul trovare i nodi più centrali, dove "centrali" può assumere diversi significati.

- **Degree centrality** : la centralità del singolo nodo viene definita come propensione ad avere legami con altri nodi; un criminale diviene centrale in quanto ha molte conoscenze.

- **Closeness centrality** : la centralità del singolo nodo viene misurata attraverso la distanza geodetica (la lunghezza del percorso minimo osservato) dagli altri nodi; il tempista diviene centrale perché riesce più facilmente a raggiungere quante più persone.
- **Betweenness centrality** : un nodo è tanto più centrale quanto più è presente negli short paths tra due generici nodi ; il criminale risulta centrale in quanto fa da ponte tra i vari livelli gerarchici.
- **Eigen vector centrality** : un nodo è tanto più centrale quanto più è collegato ad altri nodi considerati centrali; vengono identificati i criminali più alti in gerarchia.

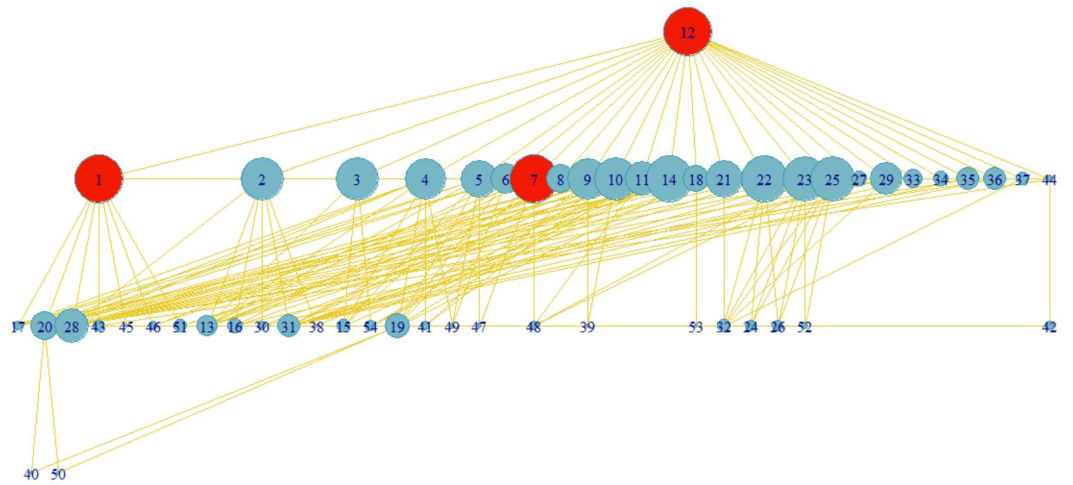
Di seguito sono riportati i network con la grandezza dei nodi proporzionale alla centralità, in rosso sarà indicata la top 3.

Notiamo come nessun nodo del terzo o quarto livello risulti centrale, ovviamente questo è dovuto al tipo di organizzazione in esame. Mentre i nodi 1, 7 e 12 risultano sempre centrali, ovviamente il 12 lo è in quanto "capo" della gerarchia, mentre gli altri due nodi potrebbero essere i suoi diretti sottoposti, infatti essi risultano più centrali del 12 in termini di betweenness centrality, mentre hanno lo stesso peso nel senso di eigen vector centrality.

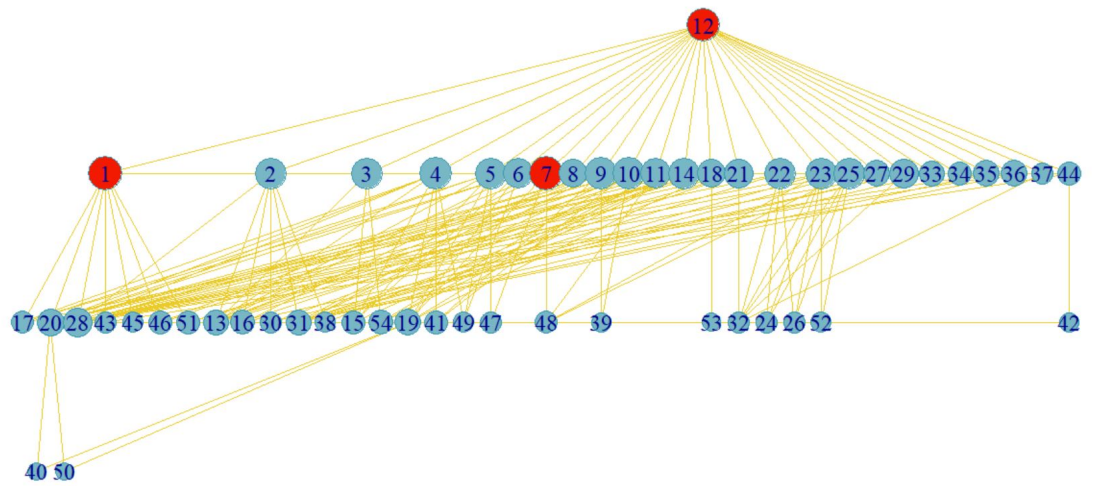
Un modo per sintetizzare queste misure è quello di creare degli indici, i quali hanno un range che va da zero ad uno. Leggendo i risultati della tabella sottostante, si conferma ancora una volta la struttura del network, infatti una Betweenness bassa indica la presenza di pochi capi, mentre la closeness rimarca la gerarchia delle associazioni a delinquere.

Tipo di centralità	CI
Degree	0.2612482
Closeness	0.1555041
Betweenness	0.003464251

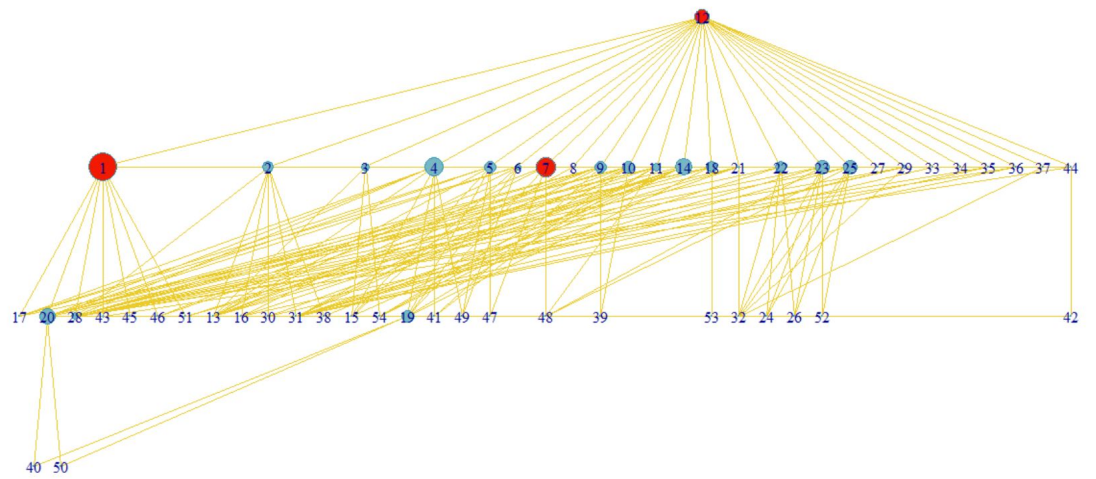
Gang network: Degree centrality



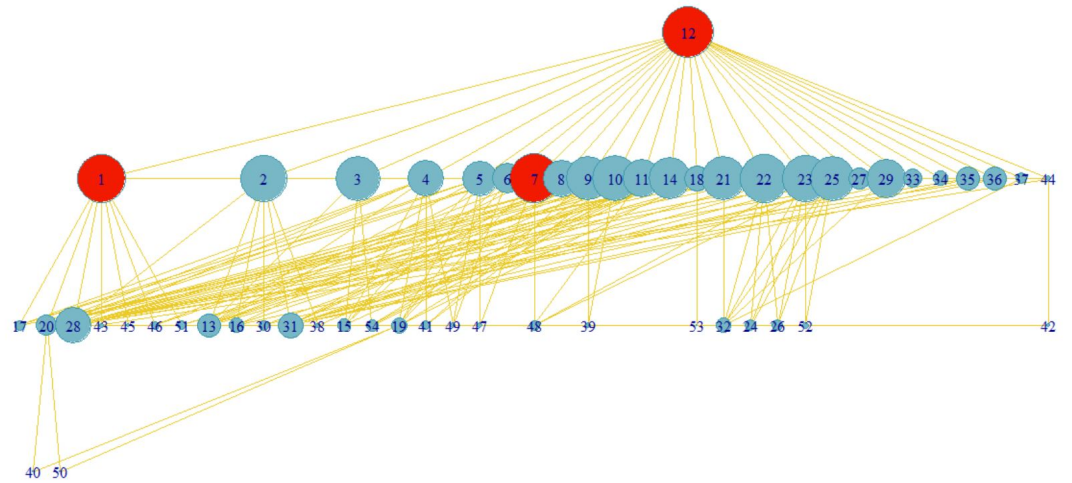
Gang network: Closeness centrality



Gang network: Betweenness centrality



Gang network: Eigenvector centrality



Chapter 2

Exponential Random Graph Model

Successivamente \mathbf{Y} indicherà la matrice di adiacenza e Y_{ij} un suo valore.

2.1 Modello Nullo: Random Graph Model

Il Random Graph Model assume che ogni nodo ha la stessa probabilità di essere osservato, indipendentemente dagli altri, ovvero:

$$Y_{ij} \sim \text{Bern}(p)$$

Questo modello risulta un punto di partenza obbligatorio per la ricerca del nostro modello migliore, in quanto ci permette di chiarire il rischio di osservare un arco.

Il modello usa un metodo di stima MLE, il termine dell'arco è il log-odds osservato (sotto l'assunzione di indipendenza), il segno negativo indica che la probabilità di osservare un arco è inferiore a 0.5.

$$\begin{aligned}
L(p) &= \prod_{i=1}^n p^{y_i} (1-p)^{(1-y_i)} \\
\ell(p) &= \log p \sum_{i=1}^n y_i + \log(1-p) \sum_{i=1}^n (1-y_i) \\
\frac{\partial \ell(p)}{\partial p} &= \frac{\sum_{i=1}^n y_i}{p} - \frac{\sum_{i=1}^n (1-y_i)}{1-p} \stackrel{\text{set}}{=} 0 \\
\sum_{i=1}^n y_i - p \sum_{i=1}^n y_i &= p \sum_{i=1}^n (1-y_i) \\
p &= \frac{1}{n} \sum_{i=1}^n y_i \\
\frac{\partial^2 \ell(p)}{\partial p^2} &= -\frac{\sum_{i=1}^n y_i}{p^2} - \frac{\sum_{i=1}^n (1-y_i)}{(1-p)^2} \\
&= -1.2649 \\
&(s.e. = 0.0638)
\end{aligned}$$

Il relativo odds é pari a 1.28, il coefficiente, risulta significativo, indi per cui la differenza tra la probabilità di osservare o meno un arco é significativa.

2.2 Modelli intermedi

2.2.1 Modello senza attributi

Nel nostro framework l'interazione con piú nodi, o specifici nodi, é centrale, in quanto definiscono il ruolo dell'unità e la sua posizione gerarchica, inoltre anche gli attributi giocano un ruolo importante. La scelta naturale del modello ricade su un modello Markoviano, in quanto assume la dipendenza delle diadi, ovvero:

$$Y_{ij} \not\perp\!\!\!\perp Y_{kl} \text{ se } \{i, j\} \cap \{k, l\} \neq \emptyset$$

Il modello di Markov viene rappresentato attraverso il grafico delle dipendenze, in particolare il legame tra $Pr(\mathbf{Y} = \mathbf{y})$ e questo tipo di grafo viene espresso nel teorema di Hammersley-Clifford. Il teorema afferma che la probabilità del modello per \mathbf{Y} é definita solo dalle cliques presenti nel suo grafo delle dipendenze.

Il modello generale risulta essere:

$$Pr(\mathbf{Y} = \mathbf{y}) = \kappa^{-1} \exp \mu L(y) + \sigma_2 S_2(y) + \sigma_3 S_3(y) + \dots + \sigma_{n-1} S_{n-1} + \tau T(y)$$

Dove S_r e T sono, rispettivamente, le r-star e i triangoli osservati.

Ovviamente per facilitare la stima dei coefficienti del modello é necessario

fare delle assunzioni di omogeneità, ovvero affermare che gli effetti delle r-str e dei triangoli siano uguali per tutti i nodi.

Come é noto il modello cosí definito ha problemi di identificabilit , indi per cui   necessario snellire il carico degli stimatori, ovvero bisogna trovare una soluzione al fatto che il modello degeneri e non arrivi a una vera e propria stima, per via dei pochi dati a disposizione.

Usando un approccio partial conditional, aggiriamo la degenery iusses e riusciamo finalmente ad avere il modello, inserendo come covariate i k-tringoli alternati e i k-star alternati.

Il modello finale scelto risulta essere quello contenente solo i k-triangoli alternati. Di seguito   riportato l'output.

```
=====
Summary of model fit
=====

Formula:   net ~ edges + gwesp(decay = 0.3, fixed = T)

Iterations: 5 out of 20

Monte Carlo MLE Results:
      Estimate Std. Error MCMC % p-value
edges      -8.4574    0.9391      0 <1e-04 ***
gwesp.fixed.0.3  4.8172    0.6843      0 <1e-04 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

      Null Deviance: 1984 on 1431 degrees of freedom
      Residual Deviance: 1391 on 1429 degrees of freedom

AIC: 1395    BIC: 1406    (Smaller is better.)
```

Il valore e il segno dei parametri conferma quanto detto in precedenza, ovvero vi   una bassa probabilit  di osservare archi, ma una tendenza positiva a costruire triangoli, in questo caso alternati, tutti tratti distintivi di network con cluster o comunque con strutture piramidali.

Insieme al modello cambia anche il metodo di stima, in questo caso viene usato l'MCMCLE, ovvero il Markov chain Monte Carlo Maximum Likelihood estimation. L'algoritmo stima un vettore di parametri attraverso un metodo di ottimizzazione per la likelihood, successivamente vengono pescate diverse matrici di adiacenza usando lo stesso vettore e infine si aggiorna il parametro, i due passi vengono iterati fino alla convergenza.

Risulta fondamentale la scelta del punto di inizio, e obbligatorio il controllo della convergenza della catena.

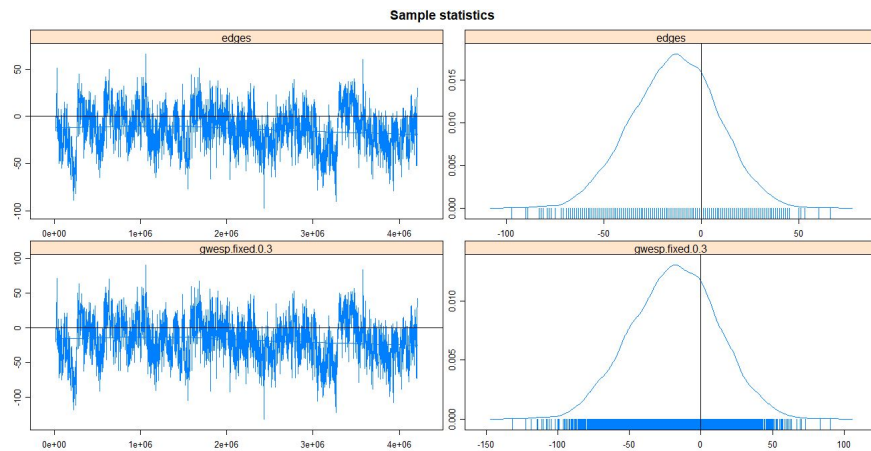


Figure 2.1: La catena sembra converga in modo corretto.

2.2.2 Modello con gli attributi

Definiamo ora l'effetto di omifilia, ovvero quello che misura quanto forte un nodo con un certo attributo x tende a essere connesso con altri nodi aventi lo stesso attributo. In questa analisi ha senso considerare solo questo effetto, che può essere visto come tendenza di connessioni all'interno di un cluster. Il modello migliore risulta essere quello contenente il luogo di nascita, gli arresti, le condanne e le età, queste tre caratteristiche sono molto legate tra di loro, infatti il luogo di nascita media tutti gli effetti.

```
=====
Summary of model fit
=====

Formula: net ~ edges + gwasp(decay = 0.3, fixed = T) + nodefactor("Birthplace",
  base = 4) + nodecov("Age") + nodecov("Arrests") + nodecov("Conviction")

Iterations: 5 out of 20

Monte Carlo MLE Results:
      Estimate Std. Error MCMC % p-value
edges          -9.43878    1.06351    0 < 1e-04 ***
gwasp.fixed.0.3    3.72046    0.65343    0 < 1e-04 ***
nodefactor.Birthplace.1 -0.02848    0.12893    0 0.825210
nodefactor.Birthplace.2 -0.28080    0.13387    0 0.036117 *
nodefactor.Birthplace.3 -0.31518    0.12907    0 0.014730 *
nodecov.Age         0.06619    0.01860    0 0.000384 ***
nodecov.Arrests     0.07950    0.01567    0 < 1e-04 ***
nodecov.Conviction  -0.14200    0.02894    0 < 1e-04 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null Deviance: 1984 on 1431 degrees of freedom
Residual Deviance: 1335 on 1423 degrees of freedom

AIC: 1351    BIC: 1394    (Smaller is better.)
```

Tutti i parametri degli attributi dicotomici positivi hanno un effetto di omofilia positivo, quelli negativi viceversa. Per Age invece, il segno positivo indica una alta attività dei nodi con un'età alta, ciò é spiegato dal fatto che chi ha più anni é più in alto nella gerarchia.

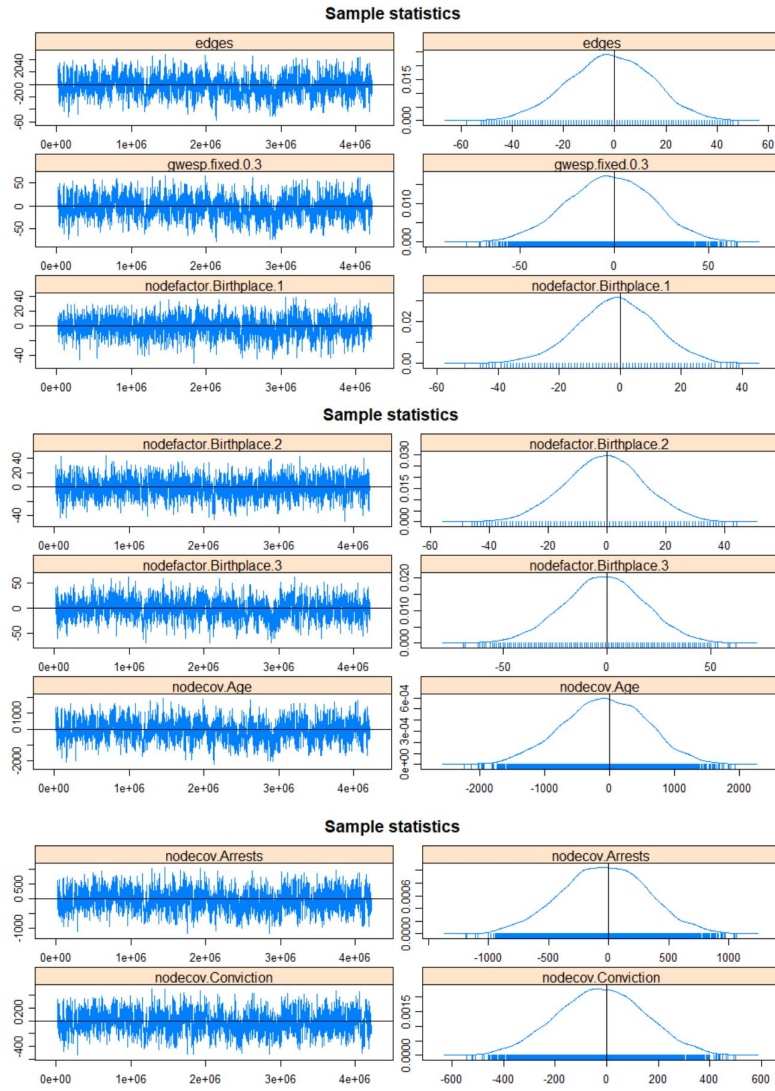


Figure 2.2: il modello sembra avere una buona convergenza

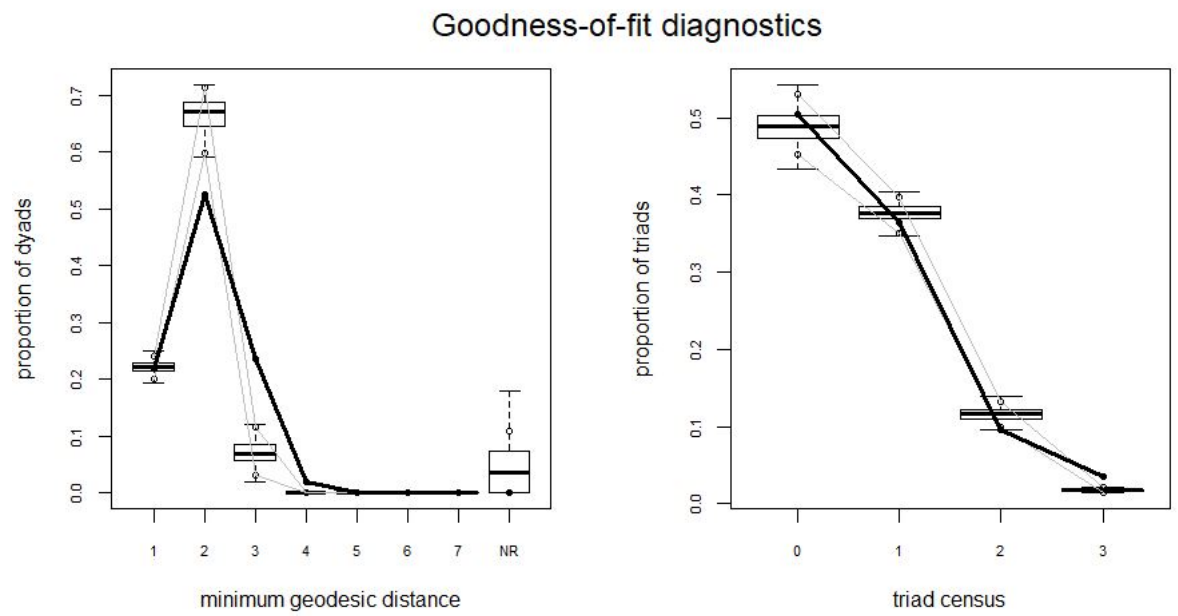


Figure 2.3: Il grafico riassume la bontà di adattamento, la quale risulta non essere buona per le distanze geodesiche, mentre è molto buona per le triadi.

Chapter 3

Modello a blocchi latenti

Con il modello a blocchi latente possiamo migliorare la comprensione del ruolo giocato dal luogo di nascita.

Anzitutto ricordiamo che questo modello ha la probabilità di osservare un legame che dipende solo dal proprio gruppo di appartenenza, in altre parole: i nodi sono indipendenti dato il gruppo di appartenenza.

Stimato come 4 il numero di cluster più adatti, osserviamo le probabilità a priori di appartenere a ciascun gruppo.

Probabilità apriori di ciascun gruppo	0.1352106	0.2608365	0.3953422	0.2086108
---------------------------------------	-----------	-----------	-----------	-----------

Mentre le probabilità a posteriori di osservare un arco tra e nei vari gruppi sono riportate nella tabella seguente.

	G1	G2	G3	G4
G1	0.572	0.478	0.047	0.402
G2	0.478	0.377	0.101	0.000
G3	0.047	0.101	0.071	0.000
G4	0.402	0.000	0.000	0.224
Numerosità	9	13	21	11

Si noti che alcune probabilità sono pari a zero, questo perché avendo solo 54 osservazioni, non è possibile effettuare una perfetta analisi.

Di seguito sono riportate le rappresentazioni grafiche dell'output della funzione *mixer()*, insieme al boxplot che spiega la relazione tra i gruppi e il luogo di nascita.

Particolarmente interessante è il confronto visivo tra il grafo colorato sotto le indicazioni del modello a blocchi latenti e quello riguardo l'attributo preso in esame. Essi risultano molto simili, ma data la bassa numerosità dei gruppi, il modello non riesce a catturare tutti i punti.

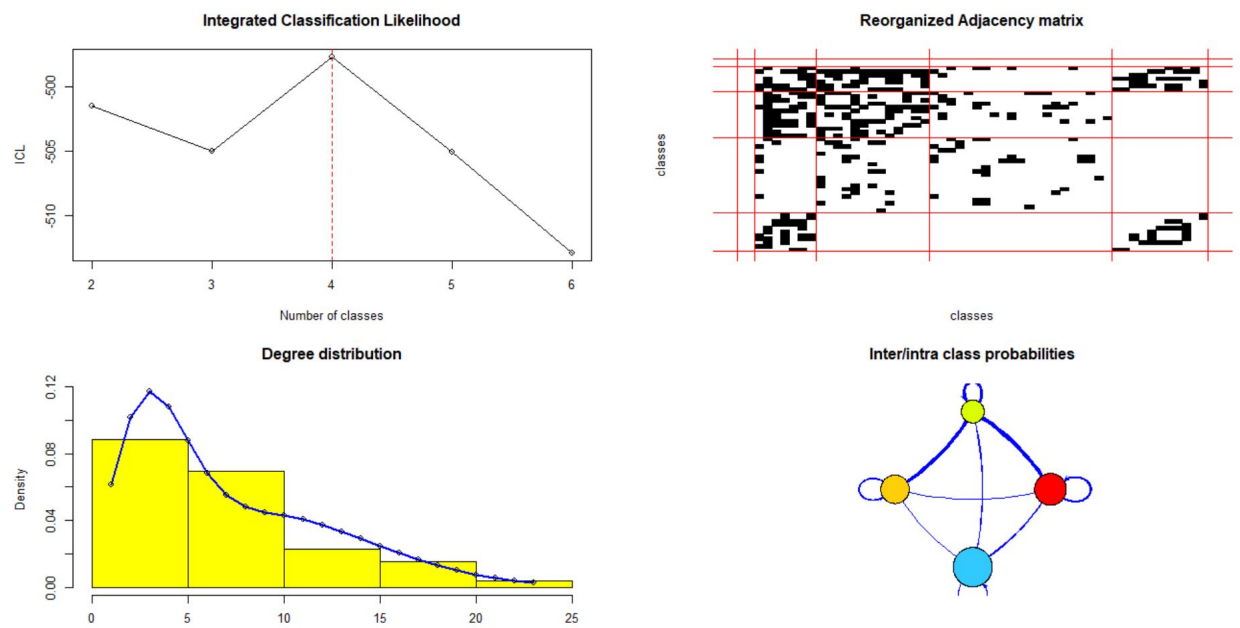


Figure 3.1: Summary grafico della funzione mixer.

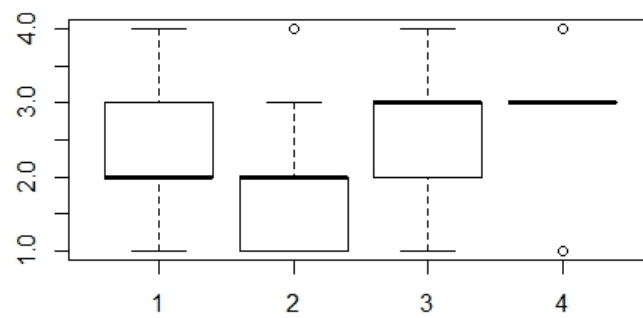


Figure 3.2: Box plot del luogo di nascita spiegato dai gruppi.

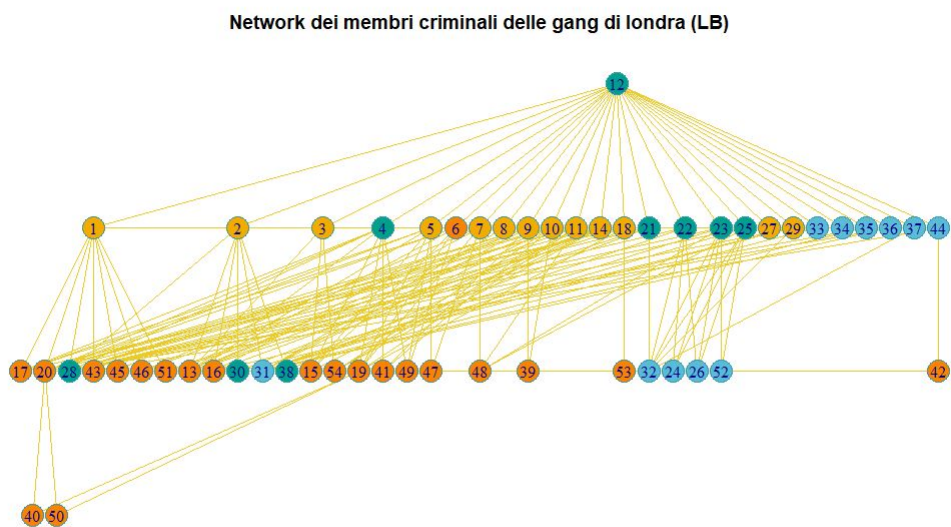


Figure 3.3: In questa rappresentazione ogni colore del nodo rappresenta un gruppo diverso

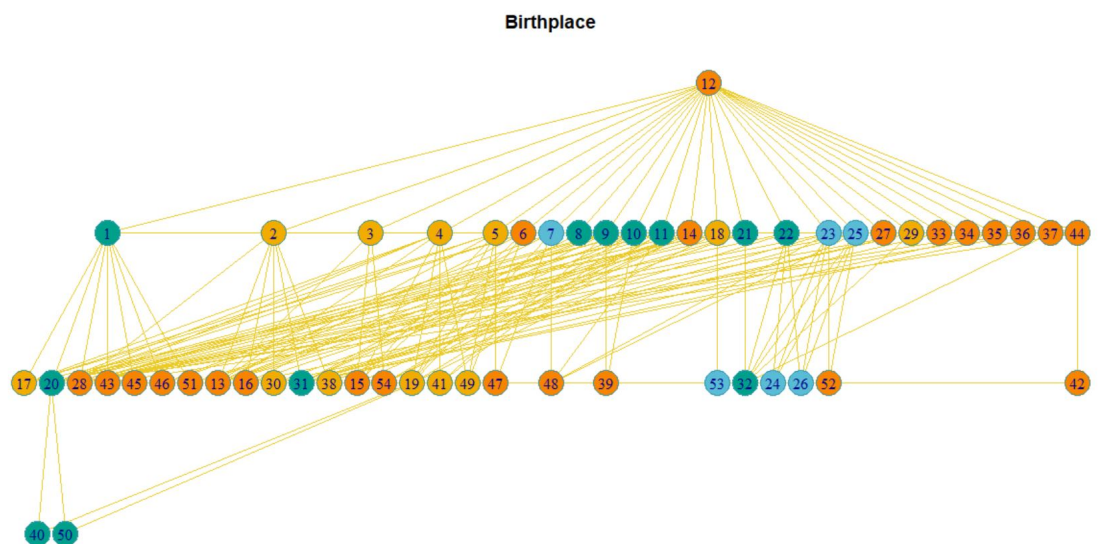


Figure 3.4: Sopra é riportato il grafo con i nodi colorati in base al luogo di nascita, il rosso indica i nati in Africa dell'ovest , il verde chi é dei Caraibi, il giallo é per i nati negli UK, l'arancione é per quelli dell'Africa dell'est

Appendice

Di seguito sono riportati gli output dei modelli che hanno contribuito a spiegare al meglio il ruolo dell'attributo *Birthplace*.

```
=====
Summary of model fit
=====

Formula: net ~ edges + gwesp(decay = 0.3, fixed = T) + nodematch("Residence") +
  nodematch("Prison") + nodefactor("Birthplace") + nodecov("Age") +
  nodecov("Arrests") + nodecov("Conviction")

Iterations: 5 out of 20

Monte Carlo MLE Results:
      Estimate Std. Error MCMC % p-value
edges      -9.58611    1.10091    0 < 1e-04 ***
gwesp.fixed.0.3  3.72718    0.67935    0 < 1e-04 ***
nodematch.Residence  0.06901    0.12953    0 0.594313
nodematch.Prison    0.07350    0.12996    0 0.571768
nodefactor.Birthplace.2 -0.26020    0.11086    0 0.019058 *
nodefactor.Birthplace.3 -0.28786    0.10029    0 0.004163 **
nodefactor.Birthplace.4  0.02665    0.12779    0 0.834834
nodecov.Age         0.06617    0.01860    0 0.000386 ***
nodecov.Arrests     0.08061    0.01600    0 < 1e-04 ***
nodecov.Conviction  -0.14345    0.02974    0 < 1e-04 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null Deviance: 1984 on 1431 degrees of freedom
Residual Deviance: 1335 on 1421 degrees of freedom
AIC: 1355 BIC: 1407 (Smaller is better.)
```

Figure 3.5: Modello completo con tutti gli attributi

```

=====
Summary of model fit
=====

Formula:  net ~ edges + gwesp(decay = 0.3, fixed = T) + nodefactor("Birthplace")

Iterations: 5 out of 20

Monte Carlo MLE Results:
      Estimate Std. Error MCMC % p-value
edges          -7.43574    0.97808      0 <1e-04 ***
gwesp.fixed.0.3    4.41752    0.68318      0 <1e-04 ***
nodefactor.Birthplace.2 -0.15939    0.09934      0  0.109
nodefactor.Birthplace.3 -0.39065    0.09798      0 <1e-04 ***
nodefactor.Birthplace.4 -0.05322    0.11835      0  0.653
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

      Null Deviance: 1984 on 1431 degrees of freedom
      Residual Deviance: 1371 on 1426 degrees of freedom

AIC: 1381    BIC: 1407    (Smaller is better.)

```

Figure 3.7: Modello con l'attributo *Birthplace*

```

=====
Summary of model fit
=====

Formula:  net ~ edges + gwesp(decay = 0.3, fixed = T) + nodecov("Age") +
      nodecov("Arrests") + nodecov("Conviction")

Iterations: 4 out of 20

Monte Carlo MLE Results:
      Estimate Std. Error MCMC % p-value
edges          -10.25468    1.02613      0 < 1e-04 ***
gwesp.fixed.0.3    3.97223    0.68973      0 < 1e-04 ***
nodecov.Age         0.06659    0.01725      0 0.000118 ***
nodecov.Arrests     0.08098    0.01569      0 < 1e-04 ***
nodecov.Conviction  -0.14380    0.02918      0 < 1e-04 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

      Null Deviance: 1984 on 1431 degrees of freedom
      Residual Deviance: 1348 on 1426 degrees of freedom

AIC: 1358    BIC: 1385    (Smaller is better.)

```

Figure 3.6: Modello completo senza l'attributo *Birthplace*

```

=====
Summary of model fit
=====

Formula:  net ~ edges + gwesp(decay = 0.3, fixed = T) + nodefactor("Birthplace") +
          nodecov("Arrests") + nodecov("Conviction")

Iterations: 5 out of 20

Monte Carlo MLE Results:

```

	Estimate	Std. Error	MCMC %	p-value
edges	-7.304798	0.946082	0	< 1e-04 ***
gwesp.fixed.0.3	3.977495	0.669269	0	< 1e-04 ***
nodefactor.Birthplace.2	-0.194073	0.105873	0	0.066999 .
nodefactor.Birthplace.3	-0.343811	0.098615	0	0.000504 ***
nodefactor.Birthplace.4	0.001219	0.125774	0	0.992267
nodecov.Arrests	0.061927	0.014487	0	< 1e-04 ***
nodecov.Conviction	-0.089606	0.024267	0	0.000230 ***

```

---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null Deviance: 1984 on 1431 degrees of freedom
Residual Deviance: 1348 on 1424 degrees of freedom

AIC: 1362    BIC: 1399    (Smaller is better.)

```

Figure 3.8: Modello completo senza l'attributo *Age*

Riferimenti

Grund, T. and Densley, J. (2015) Ethnic Homophily and Triad Closure: Mapping Internal Gang Structure Using Exponential Random Graph Models. *Journal of Contemporary Criminal Justice*, Vol. 31, Issue 3, pp. 354-370

Grund, T. and Densley, J. (2012) Ethnic Heterogeneity in the Activity and Structure of a Black Street Gang. *European Journal of Criminology*, Vol. 9, Issue 3, pp. 388-406. SOURCE: Available from Manchester.