

Università degli studi di Parma

Facoltà di Scienze Naturali, Fisiche e Matematiche

Corso di Laurea in Fisica Teorica

Sviluppi teorici e applicativi delle metriche entropiche di Rohlin

Relatore:
prof. Mario Casartelli

Correlatore:
prof.ssa Raffaella Burioni

Tesi presentata da:
Dawid Crivelli Wieslaw

Anno Accademico 2010-2011

Indice

1	Partizioni su insiemi	5
1.1	Generalità	5
1.1.1	Distanza topologica	7
1.2	Riduzione	8
1.2.1	Amplificazione	9
1.3	Possibili tipi di partizionamento	10
1.3.1	Sequenze lineari connesse	10
1.3.2	Sequenze non connesse	12
1.3.3	Sequenze su reticoli	12
1.3.4	Grafi arbitrari	14
2	Algoritmi	15
2.1	Partizioni connesse	15
2.1.1	Entropia di una partizione binaria	16
2.1.2	Distanza di Rohlin	16
2.2	Partizionamento generico	17
2.2.1	Struttura dei vicini	17
2.2.2	Partizionamento	18
2.2.3	Relabeling	20
2.2.4	Hashing	21
2.3	Operazioni tra partizioni generiche	22
2.3.1	Prodotto	22
2.3.2	Entropia	22
2.3.3	Differenza simmetrica	23
2.3.4	Intersezione	24
2.4	Riduzione	25
2.4.1	Confronto con l'intersezione	25
2.4.2	Confronto diretto	26
3	Proteine dei virus dell'influenza	27
3.1	Il virus dell'influenza	27
3.1.1	Influenza A	27
3.1.2	Vaccini	30
3.2	Clustering	31
3.2.1	Clustering gerarchico	32
3.2.2	Clustering spettrale	33
3.2.2.1	Riduzione dimensionale	34
3.2.2.2	Taglio del grafo	34
3.2.3	Numero ottimale di clusters	37
3.3	Analisi preliminare delle sequenze	40
3.4	Scelta della distanza	43
3.4.1	Distanza di Hamming	44
3.4.2	Distanza topologica	44
3.4.3	Distanza tra partizioni sconnesse	46
3.4.4	Proteine diverse, NA	50

3.5	Conclusione	52
-----	-------------------	----

Capitolo 1

Partizioni su insiemi

Definiamo in questo capitolo alcuni concetti fondamentali: le partizioni su un insieme, il significato di entropia, lo spazio delle partizioni con le operazioni binarie associate e una distanza tra partizioni. Studiamo le proprietà di questa distanza e come amplificarla. Diamo infine qualche esempio sui possibili tipi di partizioni che si possono incontrare su un insieme discreto.

1.1 Generalità

Introduciamo il formalismo e risultati generali per spazi di partizioni e metriche di Rohlin, seguendo l'approccio in [1, 2]. Sia $(\mathbf{M}, \mathcal{M}, \mu)$ uno spazio di probabilità, ovvero un insieme \mathbf{M} , una σ -algebra \mathcal{M} di sottoinsiemi di \mathbf{M} , una misura normalizzata μ su \mathcal{M} . Nei casi trattati \mathbf{M} può essere una sequenza di simboli, un reticolo bidimensionale, un grafo arbitrario.

Introducendo una relazione di equivalenza su \mathbf{M} , possiamo definire una particolare classe di sottoinsiemi. Una *partizione* di \mathbf{M} è una collezione finita $\alpha \equiv (A_1, A_2, \dots, A_N)$ di sottoinsiemi disgiunti misurabili che ricoprono \mathbf{M} , cioè $A_i \cap A_j = \emptyset$ se $i \neq j$ e $\bigcup_k A_k = \mathbf{M}$. Gli $\{A_k\}$ sono chiamati *atomi* di α e sono una rappresentano le classi di equivalenza degli elementi di \mathbf{M} . L'insieme di tutte le partizioni misurabili è denotato con $\mathcal{Z} \equiv \mathcal{Z}(\mathbf{M})$. La partizione unitaria ν consiste del singolo atomo coincidente con \mathbf{M} . È possibile introdurre un ordine parziale su \mathcal{Z} , con la relazione $\alpha \leq \beta$ quando β è un raffinamento di α : questo accade quando ogni atomo A_k è esattamente composto da atomi di β , cioè $A_k = \left\{ \bigcup_j B_j \mid B_j \in \beta \right\}$. In questo caso, si dice che α è un *fattore* di β . La partizione banale $\nu \leq \alpha, \forall \alpha$.

I termini come *unità* e *fattore* dipendono dalla definizione di uno pseudo-prodotto commutativo ed associativo, la composizione $\gamma = \alpha \vee \beta$ (o anche $\gamma = \alpha\beta$ ove non vi sia ambiguità). Il *prodotto* è la partizione meno fine di tutte le partizioni con $\gamma \geq \alpha$ e $\gamma \geq \beta$, i cui atomi sono le intersezioni non vuote di tutti gli atomi di α e β . Chiaramente il prodotto con l'unità si comporta come l'identità del prodotto, con $\alpha\nu = \alpha$ per ogni α , mentre $\alpha\eta = \alpha$ quando $\eta \leq \alpha$. Queste proprietà rendono il risultato della composizione una specie di "minimo comune multiplo". Dalla definizione segue anche che $\alpha \vee \alpha = \alpha$, ovvero il prodotto è idempotente.

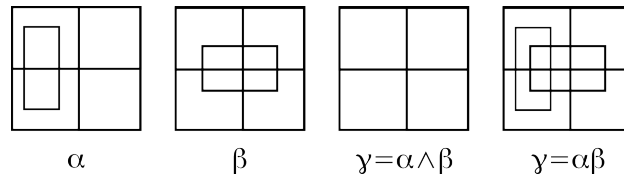


Figura 1.1.1: Esempio di prodotto e intersezione tra due partizioni

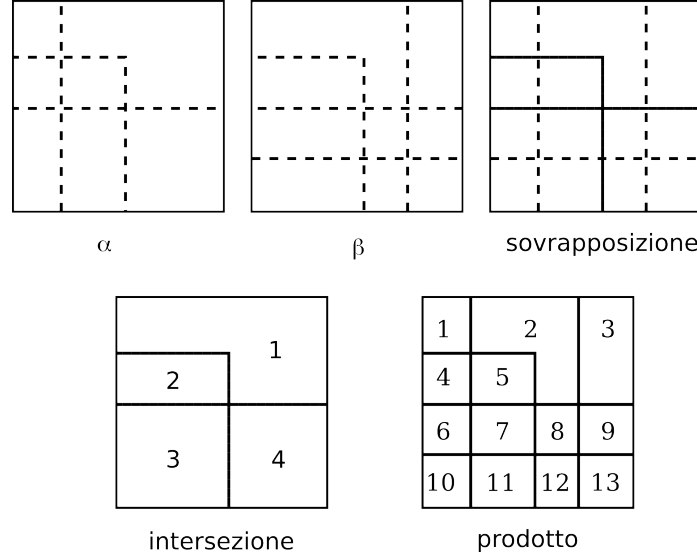


Figura 1.1.2: Intersezione e prodotto come operazioni sui bordi. Le partizioni α e β hanno bordi tratteggiati, con fase complementare, che danno l'impressione di linea continua quando sovrapposti. È evidente in questo modo il tratto di bordi in comune. Gli atomi sono stati numerati per rendere evidente la differenza tra le due operazioni.

L'operazione complementare che implementa un "massimo comune divisore" nello spazio delle partizioni è l'*intersezione* $\sigma = \alpha \wedge \beta$, definita come la partizione più fine tale che $\sigma \leq \alpha$ e $\sigma \leq \beta$. In questo caso, $\alpha \wedge v = v$ per ogni α , mentre $\alpha \wedge \beta = v$ implica che α e β sono *relativamente primi*, cioè non hanno un fattore comune.

Un altro modo di calcolare e visualizzare le operazioni è in termini di *bordi interni* di atomi della partizione, come si vede in figura 1.1.2. Il prodotto $\alpha \vee \beta$ corrisponde alla partizione avente come bordi l'unione dei bordi di α e β , mentre l'intersezione $\alpha \wedge \beta$ ha come bordi l'intersezione di quelli di α e β . Quando nell'intersezione i bordi non si chiudono, sono cancellati dal risultato e gli atomi raggruppati. Poiché la partizione banale non ha bordi tra atomi, si ricavano immediatamente le sue proprietà nel prodotto e nell'intersezione.

Una partizione può rappresentare un esperimento probabilistico con risultati disgiunti A_1, \dots, A_N , dove l'evento *atomico* A_k ha probabilità $\mu(A_k)$. Un *fattore* è quindi un sottoesperimento dell'esperimento più fine, che raggruppa diversi risultati come equivalenti: ad esempio, "pari o dispari" è un sottoesperimento con due atomi, dell'esperimento $\{1, 2, 3, 4, 5, 6\}$ del lancio di un dado.

Sullo spazio \mathcal{X} possiamo definire dei funzionali *entropia* $H : \mathcal{X} \rightarrow \mathbb{R}^+$, definiti su ogni partizione. In particolare l'*entropia di Shannon*

$$H_S(\alpha) = - \sum_{i=1}^n \mu(A_i) \ln \mu(A_i) \quad (1.1.1)$$

L'entropia di Shannon è una misura dell'informazione media ottenuta dall'esperimento. Si vede immediatamente che la partizione banale v , non codificando alcuna informazione ha entropia nulla. Se $\beta = (B_1, \dots, B_n)$ è un'altra partizione, l'entropia condizionata di α rispetto a β è

$$H_S(\alpha|\beta) = - \sum_{i=1}^n \sum_{k=1}^m \mu(A_i \cap B_k) \ln \frac{\mu(A_i \cap B_k)}{\mu(B_k)} = H_S(\alpha \vee \beta) - H_S(\beta) \quad (1.1.2)$$

dove si prende per convenzione che $x \ln x = 0$ se $x = 0$. L'entropia condizionata è l'informazione media residua ottenuta da α quando il risultato per β è noto. Si noti che l'entropia di Shannon dipende solo dalla distribuzione delle misure degli atomi, non dalla loro

natura o “forma”, che potrebbe non avere significato in spazi astratti. Le mutue relazioni tra atomi (e possibilmente le loro forme) al contrario influenzano direttamente l'entropia condizionale.

Definiamo una metrica sullo spazio delle partizioni $\mathcal{Z}(\mathbf{M})$ tramite la distanza di Rohlin d_R

$$d_R = H(\alpha|\beta) + H(\beta|\alpha)$$

che misura la complessiva non-similarità tra le partizioni α e β . È possibile dare una definizione alternativa di questa distanza, sfruttando la seconda scrittura della probabilità condizionale, riscrivendo d_R come

$$d_R = 2H(\alpha \vee \beta) - H(\alpha) - H(\beta) \quad (1.1.3)$$

La simmetria $d_R(\alpha, \beta) = d_R(\beta, \alpha)$, la positività e la condizione $d_R(\alpha, \alpha) = 0$ sono manifeste, mentre la disuguaglianza triangolare è soddisfatta se H soddisfa alle condizioni di un funzionale entropia [1]:

$$H(\alpha \vee \beta | \sigma) = H(\alpha | \beta) + H(\beta | \alpha \vee \sigma)$$

$$H(\alpha \vee \beta | \sigma) \leq H(\alpha | \sigma) + H(\beta | \sigma)$$

$$H(\alpha | \sigma) \leq H(\beta | \sigma) \quad \text{se } \alpha \subset \beta$$

$$H(\alpha | \sigma) \leq H(\alpha | \beta) \quad \text{se } \sigma \supset \beta$$

Se \mathbf{M} è finito, una *configurazione* o *stato* \mathbf{a} su \mathbf{M} è una funzione che assegna ad ogni punto $x_i \in \mathbf{M}$ un valore $a_i = f(x_i)$ nell'alfabeto \mathbb{K} . Tutte le possibili configurazioni formano uno spazio $\mathcal{C} \equiv \mathcal{C}(\mathbf{M})$. Su \mathcal{C} la distanza di Hamming è definita come

$$d_H(\mathbf{a}, \mathbf{b}) = \mathcal{N} \sum_i \rho(a_i, b_i)$$

dove $\rho(a_i, b_i)$ è una distanza su \mathbb{K} e \mathcal{N} una possibile costante di rinormalizzazione, che noi porremo uguale a 1. Da notare come questa distanza operi tra elementi diversi dalla distanza di Rohlin, pur partendo sempre da \mathbf{M} .

1.1.1 Distanza topologica

Oltre all'entropia di Shannon, definita tramite la misura μ , è possibile una definizione alternativa che non vi fa ricorso. Per questo motivo, è chiamata *entropia topologica*. Dato uno spazio compatto, come lo spazio \mathbf{M}^1 da noi considerato, il teorema di Heine-Borel[3] assicura che da ogni ricoprimento tramite insiemi aperti è possibile estrarre una sottocopertura finita. In particolare, $\exists \delta \in \mathbb{N}^+$ numero minimo di aperti con cui generare tale ricoprimento. La partizione di uno spazio è ricoprimento *minimale*², motivo per cui il numero di atomi della partizione satura la disuguaglianza del teorema, $\delta = n$. Una volta stabilito l'esistenza di δ e avendone calcolato il risultato, poniamo l'entropia topologica pari al logaritmo naturale di δ

$$H_T(\alpha) = \ln(\delta) = \ln(n) \quad (1.1.4)$$

Nei casi in cui non vi sia una naturale misura sullo spazio \mathbf{M} , la definizione topologica può essere utile per cercare una distanza che meglio renda la differenza intrinseca tra partizioni. Tuttavia, non esiste una definizione di entropia condizionale per H_T e per definire la distanza di Rohlin utilizziamo la sua definizione in termini di prodotto tra partizioni.

Si vede come H_T è un buon funzionale entropia, infatti se $\alpha < \beta$, allora $H_T(\alpha) < H_T(\beta)$, in quanto una partizione strettamente più fine ha banalmente un numero maggiore di atomi. Da questa condizione si ha che per il prodotto

$$\alpha \vee \beta \geq \alpha \implies H_T(\alpha \vee \beta) \geq H_T(\alpha)$$

¹Addirittura \mathbf{M} nel nostro caso è discreto, per cui la compattezza è equivalente all'essere finito.

²Nulla vieta di avere altri ricoprimenti perfettamente ridondanti, ma sempre con un numero finito di aperti

per α, β generici e la distanza di Rohlin topologica, definita usando il funzionale H_T nell'equazione 1.1.3 risulta definita positiva. Notiamo che la distanza è ben definita solo quando si considera partizioni poichè la distanza tra ricoprimenti non minimali può non essere 0 anche tra ricoprimenti diversi – ha senso solo considerare la restrizione alla classe di equivalenza del ricoprimento minimale corrispondente.

1.2 Riduzione

L'essenziale dissimilarità tra due partizioni potrebbe essere confusa ed indebolita dalla presenza di un fattore comune ampio, come ad esempio accade se gli atomi della partizione hanno lunghezza media molto breve, nel qual caso la maggioranza dei confini risulta essere la stessa. Si cerca quindi di eliminare fattori comuni il più possibile, con una *riduzione* che ci si aspetta aumenti la distanza relativa. Tuttavia questa operazione, analoga alla riduzione in minimi termini per frazioni, non è unicamente definita, in quanto le partizioni, a differenza degli interi, non ammettono una univoca fattorizzazione in fattori primi. Il ruolo dei fattori primi (ovvero fattori irriducibili) è giocato dalle sottopartizioni *dicotomiche*, che sono tuttavia ancora estremamente ridondanti ($2^{n-1} - 1$ per partizioni con n atomi).

A partire dalla partizione $\alpha \equiv (A_1, \dots, A_n)$, definiamo quindi una famiglia ristretta $\mathbf{E}(\alpha)$ di *fattori dicotomici elementari* $\tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_n$ con le seguenti caratteristiche:

1. $\mathbf{E}(\alpha)$ deve essere ben definita per ogni $\alpha \in \mathcal{X}$
2. $\mathbf{E}(\alpha)$ non deve contenere più di n (il numero di atomi in α) fattori elementari
3. $\bigvee_{k=1}^n \tilde{\alpha}_k = \alpha$

Una scelta universale consiste nel prendere come fattori dicotomici $\tilde{\alpha}_k \equiv (A_k, A_k^c)$ le partizioni formate dai singoli atomi e dai loro complementi in \mathbf{M} . Fattori di questo tipo sono chiamati *universali semplici*.

Una volta che per due partizioni α e β le famiglie di fattori dicotomici $\mathbf{E}(\alpha)$ e $\mathbf{E}(\beta)$ sono state definite, abbiamo diversi possibili processi di riduzione.

Definizione 1. Riduzione tramite partizione comune

1. Si definisce il massimo fattore comune $\sigma = \alpha \wedge \beta$
2. Si tralasciano da $\mathbf{E}(\alpha)$ e $\mathbf{E}(\beta)$ i fattori che **non** sono relativamente primi con σ , e indichiamo i fattori rimanenti come $\hat{\alpha}_k$ e $\hat{\beta}_k$ rispettivamente. Questo vuol dire che $\hat{\alpha}_k \wedge \sigma = \hat{\beta}_j \wedge \sigma = \nu$.

Definizione 2. Riduzione con eliminazione fattori in comune

1. Si tralasciano da $\mathbf{E}(\alpha)$ e $\mathbf{E}(\beta)$ i fattori che compaiono in entrambe le partizioni. Se indichiamo i fattori rimanenti come $\hat{\alpha}_k$ e $\hat{\beta}_k$ rispettivamente, questo vuol dire che $\forall A_k, \nexists B_j \in \beta | A_k = B_j$ e viceversa.

Definizione 3. Riduzione con eliminazione fattori simili

1. Si tralasciano da $\mathbf{E}(\alpha)$ e $\mathbf{E}(\beta)$ i fattori che hanno un corrispondente “simile” nell'altra partizione. Questo vuol dire che scartiamo da $\hat{\alpha}$ il fattore

$$A_k \text{ se } \exists B_j \in \beta, \text{ tale che } \mu(A_k \triangle B_j) \leq \epsilon$$

e viceversa da $\hat{\beta}$ il fattore

$$B_k \text{ se } \exists A_j \in \alpha, \text{ tale che } \mu(B_k \triangle A_j) \leq \epsilon$$

Il simbolo $A_k \triangle B_j$ indica la differenza simmetrica tra i due atomi, ovvero i siti che appartengono ad un atomo ma non all'altro.

Alla fine, per tutti i tipi di riduzione, definiamo le partizioni ridotte come $\hat{\alpha} = \vee_k \hat{\alpha}_k$ e $\hat{\beta} = \vee_k \hat{\beta}_k$, ovvero il prodotto dei fattori dicotomici “sopravvissuti”. Nel capitolo sugli algoritmi presenteremo metodi ottimali per il calcolo dei fattori dicotomici per ogni criterio presentato, che presentano una notevole complessità se eseguiti nel modo naïve.

Per il resto della sezione concentreremo la nostra attenzione sulla riduzione tramite fattore comune massimo. Si motiva la scelta del confronto con il fattore comune poichè vi sono casi in cui le partizioni non hanno atomi in comune, ma ciononostante si ha che $\sigma \neq \nu$. Questo accade, per esempio, quando $\alpha < \beta$ strettamente e non vi sono fattori comuni elementari. In questo caso allora $\sigma = \alpha$ e $\hat{\alpha} = \nu$ con questo metodo di riduzione, mentre $\hat{\alpha} = \alpha$ tralasciando i fattori comuni. Può capitare inoltre che anche se le partizioni sono già ridotte, non sono prime tra di loro. In particolare, $H(\sigma)$ indica la similarità tra le partizioni.

Osservazione La mappa simmetrica definita dalla distanza tra sequenze ridotte

$$\begin{aligned} \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ \alpha \otimes \beta &\mapsto d_R(\hat{\alpha}, \hat{\beta}) \end{aligned}$$

non definisce una distanza su \mathcal{X} . Infatti il passaggio $\alpha \rightarrow \hat{\alpha}$ dipende strettamente da β e viceversa, non è quindi possibile scrivere

$$d_R(\hat{\alpha}, \hat{\beta}) = \hat{d}_R(\alpha, \beta)$$

con \hat{d}_R una distanza.

1.2.1 Amplificazione

Il processo di riduzione porta alla definizione di partizioni con complessità possibilmente inferiore, ovvero $H(\hat{\alpha}) \leq H(\alpha)$. Questo va nel verso opposto quando si considera l'effetto sulla distanza, che invece aumenta.

Il rapporto di *amplificazione* R misura quanto la riduzione ha messo in risalto la differenza tra partizioni. Ne dimostriamo la caratteristica fondamentale:

$$R = \frac{d_R(\hat{\alpha}, \hat{\beta})}{d_R(\alpha, \beta)} \geq 1$$

Proposizione. $d_R(\hat{\alpha}, \hat{\beta}) \geq d_R(\alpha, \beta)$

Dimostrazione. Ricordando che $\sigma = \alpha \wedge \beta$ possiamo scrivere $\alpha = \sigma \hat{\alpha}$ e $\beta = \sigma \hat{\beta}$: infatti σ contiene tutti i fattori tralasciati durante la riduzione. Utilizzando ora l'idempotenza del prodotto, $\sigma = \sigma \sigma$, possiamo riscrivere la tesi utilizzando l'equazione (1.1.3)

$$2H(\sigma \hat{\alpha} \hat{\beta}) - H(\sigma \hat{\alpha}) - H(\sigma \hat{\beta}) \leq 2H(\hat{\alpha} \hat{\beta}) - H(\hat{\alpha}) - H(\hat{\beta})$$

scambiando l'ordine dei termini si ottiene

$$2H(\sigma \hat{\alpha} \hat{\beta}) - 2H(\hat{\alpha} \hat{\beta}) \leq H(\sigma \hat{\alpha}) - H(\hat{\alpha}) + H(\sigma \hat{\beta}) - H(\hat{\beta})$$

e sfruttando la formula (1.1.2) per l'entropia condizionata, la tesi si riduce a

$$2H(\sigma | \hat{\alpha} \hat{\beta}) \leq H(\sigma | \hat{\alpha}) + H(\sigma | \hat{\beta})$$

ma questo è chiaramente vero, in quanto

$$H(\sigma | \hat{\alpha} \hat{\beta}) \leq H(\sigma | \hat{\alpha}) \quad \text{e} \quad H(\sigma | \hat{\alpha} \hat{\beta}) \leq H(\sigma | \hat{\beta})$$

per le proprietà dell'entropia, poichè il termine condizionante è sicuramente maggiore, ovvero

$$\hat{\alpha} \hat{\beta} \geq \hat{\alpha} \quad \text{e} \quad \hat{\alpha} \hat{\beta} \geq \hat{\beta}$$

□

Da notare che la dimostrazione vale per la distanza di Rohlin definita sia tramite H_S che H_T .

Risulta importante la scelta della famiglia di fattori dicotomici $E(\alpha)$, che è dettata dalla topologia e geometria dello spazio delle configurazioni. La scelta della famiglia di fattori dicotomici universali semplici è sempre possibile, poichè la determinazione di A_k^c a partire da A_k è un'operazione ben definita in qualunque spazio topologico. Prendere come fattori elementari la parte interna di contorni di cluster ad esempio, richiede un concetto di orientabilità e la possibilità di definire contorni, ovvero insiemi con codimensione 1 su varietà – mentre vorremmo estendere l'analisi anche a grafi generici, privi di strutture geometriche predefinite. Già nel caso lineare è possibile prendere fattori dicotomici diversi e algoritmicamente più performanti, a patto di restringere lo studio alle partizioni con atomi formati da cluster connessi.

1.3 Possibili tipi di partizionamento

L'applicabilità dei metodi discussi è assolutamente generica, estendibile a qualunque spazio di probabilità finito si voglia considerare, vediamo dunque di dare esempi dei possibili spazi \mathbf{M} su cui abbiamo lavorato, con i relativi fattori dicotomici e conseguenze computazionali.

Essendo lo studio svolto su calcolatore, lo spazio \mathbf{M} e la sua σ -algebra sono finiti e discreti. I siti appartenenti ad \mathbf{M} possono essere sempre numerati ordinati in modo opportuno x_i , $i \in \{1, \dots, L\}$, dove L è il numero totale di siti, che si tratti di un reticolo o di un grafo. Poichè la partizione induce una relazione di equivalenza sull'insieme \mathbf{M} , indichiamo che due siti sono equivalenti se appartengono allo stesso atomo

$$i \sim j \iff \exists A_k \text{ tale che } i \in A_k, j \in A_k$$

Per partizionare i siti in atomi disgiunti, richiediamo che la *configurazione* (o *stato*) \mathcal{C} , associ ad ogni sito una lettera dell'alfabeto \mathbb{K} , considerato finito, $|\mathbb{K}| < \infty$. Il simbolo associato al sito i -esimo, non essendovi ambiguità, sarà d'ora in poi indicato con $f(i)$. Se lo stato del sistema è descritto con variabili continue (o vi è un numero enorme di possibili lettere nell'alfabeto, si pensi ad una variabile a 64 bit rappresentante un numero reale), si può sempre ridurre l'alfabeto ad un insieme $\{k_j\}$ discretizzando i valori della configurazione con il criterio:

$$f(i) := \bar{k} \quad \text{tale che} \quad |f(i - \bar{k})| = \min_{k_j \in \mathbb{K}} |f(i) - k_j|$$

oppure mettendo nello stesso atomo siti vicini che hanno valori distanti meno di ϵ :

$$i \sim j \iff |f(i) - f(j)| \leq \epsilon$$

1.3.1 Sequenze lineari connesse

Consideriamo sequenze lunghe L , provenienti da due casi:

- Problemi di meccanica statistica, in cui la configurazione è una variabile aleatoria, generata algoritmicamente da una catena di Markov in base al modello di Ising monodimensionale, nel qual caso l'alfabeto corrisponde a $\{-1, +1\}$.
- Sequenze di origine biologica, in particolare sequenze di amminoacidi (proteine), in cui $|\mathbb{K}| = 25$ (dettagli a pagina 41)

La distanza di Hamming d_H è sensibile solo a variazioni puntuali dei valori in \mathcal{C} . I siti che compongono la sequenza non si influenzano, una variazione su un sito può solo variare di $\{0, +1\}$ la distanza totale.

Ad ogni configurazione possiamo associare una partizione in \mathcal{Z} , in cui gli atomi sono formati dai cluster, cioè sottoinsiemi connessi di \mathbf{M} a valori omogenei in \mathbb{K} . Questo stabilisce una mappa $\Phi: \mathcal{C} \rightarrow \mathcal{Z}$ da ogni configurazione ad una partizione corrispondente,

rendendo possibile il confronto tra $d_H(\mathbf{a}, \mathbf{b})$ in \mathcal{C} e $d_R(\alpha, \beta)$ in \mathcal{Z} , dove $\alpha = \Phi(\mathbf{a})$ e $\beta = \Phi(\mathbf{b})$. La relazione è chiaramente del tipo multi-a-uno, infatti assegnando ad un segmento di lettere omogeneo in \mathcal{C} un diverso simbolo, non cambia la partizione corrispondente. In simboli, esprimiamo la relazione come

$$i \sim j \iff j = i + 1, f(i) = f(j)$$

È evidente quindi che variazioni locali, ad esempio il cambiamento di un singolo simbolo, possono non modificare affatto la partizione

$$\{\dots, T, T, C, A, A, \dots\} \stackrel{\Phi}{\equiv} \{\dots, T, T, M, A, A, \dots\}$$

che presenta sì una perdita di informazione, ma permette quindi anche di filtrare molto “rumore” e si è dimostrata una ottima scelta nel caso biologico e una necessità nello studio di sequenze di Ising – la hamiltoniana del sistema non distingue i valori associati ai singoli siti, ma solo differenze tra siti vicini. La suddivisione in partizioni si comporta nello stesso modo, estraendo quindi solo informazioni fisicamente rilevanti.

Esempio. Supponiamo di aver partizionato una stringa come $\{A, A, A, B, B, C, C, D, D, D\}$ o $\{+, +, +, -, -, +, +, -, -, -\}$, aventi la stessa partizione

$$\alpha = \{(1, 2, 3), (4, 5), (6, 7), (8, 9, 10)\}$$

il calcolo esplicito dell'entropia è il seguente:

$$\begin{array}{ll} A_1 = (1, 2, 3) & \mu(A_1) = \frac{3}{10} \\ A_2 = (4, 5) & \mu(A_2) = \frac{2}{10} \\ A_3 = (6, 7) & \mu(A_3) = \frac{2}{10} \\ A_4 = (8, 9, 10) & \mu(A_4) = \frac{3}{10} \end{array}$$

Ora, $H_S(\alpha) = -2(0.3)\ln(0.3) - 2(0.2)\ln(0.2) \simeq 1.36$, mentre $H_T(\alpha) = \ln(4) \simeq 1.38$. Come si vede i valori sono abbastanza simili.

Fattori dicotomici nel caso lineare connesso

Per una partizione $\alpha \equiv (A_1, \dots, A_n)$, definiamo i fattori dicotomici $\tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_n \in \mathbf{E}(\alpha)$ nel modo seguente:

$$\tilde{\alpha}_k = \left\{ \bigcup_{j=1}^k A_j, \bigcup_{j=k+1}^n A_j \right\}$$

dove ricordiamo che n è il numero di atomi. In termini di siti quindi, il fattore dicotomico k -esimo corrisponde alla partizione con tutti i siti corrispondenti ai primi k atomi presi insieme, ed il complementare corrispondente ai siti a partire dall'atomo $(k+1)$ -esimo fino ad L , come mostrato in figura 1.3.1.

Con questa particolare scelta, resa possibile dalla topologia connessa e ordinata, poiché abbiamo l'ordine ereditato da \mathbb{N} , il processo di riduzione è estremamente semplificato ed efficiente. Partendo da una partizione α , la partizione ridotta $\hat{\alpha}$ è quella in cui sono stati rimossi i bordi in comune tra α e σ , come si vede in figura 1.3.2

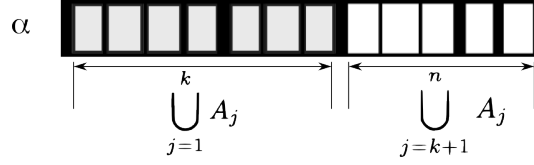
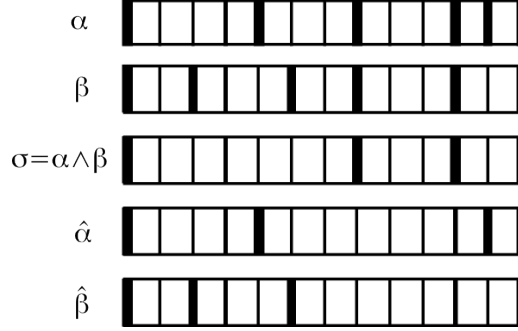


Figura 1.3.1: Fattori dicotomici lineari

Figura 1.3.2: Riduzione nel caso di sequenze lineari. La riga spessa indica la separazione tra atomi. Si nota come i bordi della partizione comune siano quelli comuni sia ad α che a β , mentre nelle partizioni ridotte non compaiono più i bordi esistenti anche in σ

1.3.2 Sequenze non connesse

Lo studio delle sequenze biologiche, in cui vi è una notevole frammentazione, dovuta alla bassa probabilità di avere molti simboli consecutivi uguali, ha motivato la ricerca di diversi criteri di partizionamento. Il più semplice è considerare come appartenenti allo stesso atomo siti con lo stesso simbolo, non adiacenti ma con la possibilità di saltare al massimo n_s siti in avanti

$$i \sim j \iff |j - i| \leq n_s, f(i) = f(j)$$

In questo caso, la topologia nello spazio \mathcal{Z} non è più connessa, per cui l'utilizzo dei fattori dicotomici lineari non è possibile. Un'alternativa è quella di prendere per l'atomo $A_k = \{n_{k_1}, n_{k_2}, \dots, n_{k_n}\}$ una partizione dicotomica costituita da tutti i siti tra il primo e l'ultimo appartenente a A_k e dal complementare, sconnesso, A_k^c

$$\hat{a}_k = \left\{ \bigcup_{i=n_{k_1}}^{n_{k_n}} n_i, \left(\bigcup_{i < n_{k_1}} n_i \right) \cup \left(\bigcup_{i > n_{k_n}} n_i \right) \right\}$$

Dati gli scarsi vantaggi di questo approccio, si è scelto tuttavia anche in questo caso di utilizzare i *fattori universali semplici*. Una volta venuta a meno la connessione degli atomi, questo tipo di partizionamento presenta le complessità e caratteristiche di un criterio arbitrario su grafo generico – si veda il prossimo paragrafo per le illustrazioni delle operazioni tra partizioni di questo tipo.

1.3.3 Sequenze su reticoli

In questo caso, si considera le partizioni in cui fanno parte dello stesso atomo siti con lo stesso simbolo, ma con il vincolo che devono essere *primi vicini* sul reticolo

$$i \sim j \iff d(i, j) = 1, f(i) = f(j)$$

dove la distanza indica il numero di passi sul reticolo. Possiamo tuttavia generalizzare arbitrariamente la condizione di vicinanza, per includere secondi vicini, siti con lo stesso simbolo arbitrariamente posti, ecc.

1	4	7
2	5	8
3	6	9

Figura 1.3.3: Mappa in memoria di un reticolo bidimensionale. In grigio e bianco sono evidenziati i valori dalla configurazione, mentre i numeri nell'angolo rappresentano l'ordine dei siti.

1	1	3
2	4	3
2	2	4

1	1	$\bar{1}$
$\bar{1}$	$\bar{1}$	$\bar{1}$
$\bar{1}$	$\bar{1}$	$\bar{1}$

$\bar{3}$	$\bar{3}$	3
$\bar{3}$	$\bar{3}$	3
$\bar{3}$	$\bar{3}$	$\bar{3}$

$\bar{2}$	$\bar{2}$	$\bar{2}$
2	$\bar{2}$	$\bar{2}$
2	2	$\bar{2}$

$\bar{4}$	$\bar{4}$	$\bar{4}$
$\bar{4}$	4	$\bar{4}$
$\bar{4}$	$\bar{4}$	4

Figura 1.3.4: Partizione con 4 atomi e i 4 fattori dicotomici universali semplici corrispondenti. I numeri indicano l'etichetta corrispondente ad ogni atomo, scelto arbitrariamente. Gli atomi sono del tipo nonconnesso, in cui le etichette indicano l'appartenenza; similmente i fattori dicotomici contengono gli atomi sconnessi ed il loro complementare indicato dal label tildato.

Illustriamo le possibili operazioni con un reticolo bidimensionale, che tuttavia generalizza immediatamente al grafo arbitrario. In particolare numeriamo anche in questo caso i siti, cosa necessaria per mappare nella memoria di un computer la configurazione. La mappatura in memoria (figura 1.3.3) comporta necessariamente la disposizione di siti "vicini" sul reticolo su posizioni in cui non risultano più contigue in memoria³. Nel caso della figura la partizione corrispondente è

$$A = \{(1, 4, 5, 6, 8, 9), (2, 4), (7)\}$$

evidentemente non-connessa. In questo caso i fattori dicotomici anche risultano sconnessi, come nell'esempio di un reticolo 3x3 in figura 1.3.4, in cui si è preso come atomi i siti con lo stesso valore indipendentemente dalla loro posizione.

Nel caso in cui gli atomi sconnessi il prodotto tra partizioni rimane banale da implementare mentre l'intersezione è estremamente complicata – un'implementazione ottimale è data nel capitolo sugli algoritmi.

³Ricordiamo la convenzione "per colonne" di rappresentazione di una matrice bidimensionale nella maggior parte dei linguaggi di programmazione, a partire dal Fortran in poi

1.3.4 Grafi arbitrari

L'insieme dei siti in questo caso è sempre ordinabile con $i \in \{1, \dots, L\}$. Definiamo la *matrice di adiacenza* A_{ij}

$$A_{ij} = \begin{cases} 1 & \text{se i siti sono connessi} \\ 0 & \text{altrimenti} \end{cases}$$

che ci permette di dare la relazione di equivalenza come

$$i \sim j \iff A(i, j) = 1, f(i) = f(j)$$

In questo modo ci si slega completamente da qualunque nozione di vicinanza geometrica. In linea di principio, è possibile rimuovere anche la condizione $f(i) = f(j)$, e lasciare alla matrice di adiacenza l'implementazione della relazione di equivalenza. Non vi sono condizioni da porre su A_{ij} , qualunque scelta è accettabile e va a modificare il partizionamento, senza che questo risulti mai incoerente.

Capitolo 2

Algoritmi

Dopo aver visto i possibili tipi di partizione, mostriamo come calcolare qualunque quantità richiesta in modo ottimale. Definiamo prima le operazioni limitate a partizioni con atomi connessi, sia per la loro minore complessità, sia perché nel caso generico si sfrutta la riduzione al caso più semplice già trattato.

2.1 Partizioni connesse

Nel caso monodimensionale è sufficiente identificare il *bordo sinistro* di ogni atomo, per ricostruire completamente la partizione. Partendo da una sequenza in un array (o segmento unidimensionale finito) $S[i]$ lungo L , definiamo l'array binario della partizione $B[i]$, lungo L , con indici che vanno da 0 a $L-1$

$$B[i] = \begin{cases} 1 & \text{il sito } i \text{ identifica un nuovo atomo} \\ 0 & \text{altrimenti} \end{cases}$$

algoritmicamente si procede in due step:

1. il primo sito sicuramente inizia un atomo, quindi $B[0] = 1$;
2. iterando su i , poniamo $B[i] = (S[i] \neq S[i-1])$; infatti un nuovo atomo corrisponde ad un valore della sequenza *diverso* dal precedente

Ogni partizione espressa in termini binari ha sempre 1 come primo elemento.

Facendo un esempio, ad una stringa possiamo semplicemente far corrispondere l'array binario corrispondente

```
S=AAAABBBAAAACCCDDDD
B=100010010001001000
```

Avendo un array che indica esattamente dove si trovano i bordi dei vari atomi, le operazioni di intersezione e prodotto sono immediate. Il prodotto corrisponde alla partizione che contiene i bordi di entrambe le partizioni (OR binario); l'intersezione ha un bordo quando compare contemporaneamente nelle due partizioni (AND binario); il prodotto ridotto contiene invece ogni bordo di una e una sola delle partizioni (XOR binario) – se il bordo appartiene ad entrambe, appartiene al fattore comune, che è stato ridotto.

In tabella 2.1 sono mostrate tutte le sequenze generabili a partire da due partizioni binarie A e B , con le rispettive entropie, che adesso mostriamo come calcolare.

Vettore	Risultato	Codice C	H_S	H_T
α	101100011001101101	$A[] = \{1, 0, 1, 1, \dots\};$	2.16	2.30
β	100010010011001000	$B[] = \{1, 0, 0, 0, \dots\};$	1.73	1.79
$\sigma = \alpha \wedge \beta$	100000010001001000	$C[i] = A[i] \ \& \ B[i];$	1.33	1.39
$\alpha \vee \beta$	101110011011101101	$P[i] = A[i] \ \ B[i];$	2.40	2.49
$\hat{\alpha} \vee \hat{\beta}$	101110001010100101	$Prid[i] = A[i] \ \wedge \ B[i];$	2.09	2.19
$\hat{\alpha}$	101100001000100101	$Arid[i] = A[i] \ \wedge \ C[i];$	1.80	1.94
$\hat{\beta}$	100010000010000000	$Brid[i] = B[i] \ \wedge \ C[i];$	1.06	1.09

Tabella 2.1: Rappresentazione delle operazioni a partire da partizioni binarie. In rosso gli elementi riconducibili a partizioni di α , in blu quelli dovuti a β , in verde gli elementi di σ . Nel prodotto si vede l'effetto sovrapposto di blu e rosso. Nelle partizioni ridotte in verde sono stati evidenziati gli zeri dovuti alla riduzione, ovvero in corrispondenza agli atomi (1 verdi) di σ .

2.1.1 Entropia di una partizione binaria

Il calcolo più semplice è quello dell'entropia topologica, H_T . Avendo un vettore che indica esattamente l'inizio di ogni atomo, basta contare il numero di 1 presenti nel vettore e farne il logaritmo¹.

Per l'entropia di Shannon H_S la cosa è leggermente più complicata²:

1. si trovano tutte le posizioni i_k degli 1 nel vettore dato
2. si fa la differenza degli interi così trovati, $\mu_k = i_k - i_{k-1}$
3. l'insieme $\{\mu_k\}$ rappresenta le lunghezze degli atomi, ovvero gli intervalli tra 1 e l'ultimo 0 (se esiste) dello stesso atomo

Per l'ultimo atomo ci vuole un trattamento speciale, $\mu_{end} = L - i_{end}$, ovvero si calcola la differenza tra la lunghezza della sequenza e l'ultimo 1 nel vettore. Il calcolo dell'entropia del vettore $\{\mu_m\} = \{\mu_k\} \cup \mu_{end}$ è molto semplice:

$$L = \sum_m \mu_m$$

$$H = - \frac{\sum_m \mu_m \ln \mu_m}{L} + \ln L$$

la formula utilizzata permette di non dover dividere ogni μ_k per L , e di eseguire un'unica divisione alla fine. L'utilizzo delle misure nonnormalizzate degli atomi consente anche di poter precalcolare tutti i logaritmi necessari, in quanto $0 \leq \mu_k \leq L$, migliorando la velocità di esecuzione della routine del 600%.

L'algoritmo presenta complessità $\mathcal{O}(L)$ in quanto si scorre l'array completo, anche se una volta sola: chiaramente non si può far di meglio.

2.1.2 Distanza di Rohlin

Una volta che si è in grado di fare il prodotto, la riduzione e il calcolo dell'entropia, si può immediatamente calcolare la distanza di Rohlin. Per questo motivo è necessario calcolare anche le partizioni ridotte, che sono superflue per il calcolo del prodotto ridotto, ma nella distanza entrano in gioco le entropie dei fattori.

In tabella 2.2 mostriamo le 4 possibili distanze date le due partizioni binarie, sfruttando solamente i numeri calcolati nell'esempio precedente.

¹in Matlab è immediato: $H_{top} = \log(\text{sum}(B));$

²In codice Matlab, si definisce prima $H = \mathcal{O}(a) - \text{sum}(a .* \log(a)) / \text{sum}(a) + \log(\text{sum}(a));$
Avendo il vettore binario A, l'entropia cercata è: $H(\text{diff}(\text{find}([A, 1])))$;

Distanza	Formula	Valore	R
Nonridotta, Shannon	$2H_S(\alpha \vee \beta) - H_S(\alpha) - H_S(\beta)$	0.93	-
Ridotta, Shannon	$2H_S(\hat{\alpha}\hat{\beta}) - H_S(\hat{\alpha}) - H_S(\hat{\beta})$	1.33	1.46
Nonridotta, topologica	$2H_T(\alpha \vee \beta) - H_T(\alpha) - H_T(\beta)$	0.87	-
Ridotta, topologica	$2H_T(\hat{\alpha}\hat{\beta}) - H_T(\hat{\alpha}) - H_T(\hat{\beta})$	1.35	1.53

Tabella 2.2: Tutte le possibili distanze ricavate a partire dalla tabella 2.1, con relativo fattore di amplificazione

2.2 Partizionamento generico

Nel caso generico non basta indicare l'inizio di ogni atomo, poichè questi possono avere forma (e buchi) di ogni tipo. È necessario etichettare ogni sito con un numero naturale, indicando con $\lambda_\alpha(i)$ l'etichetta corrispondente al sito i -esimo nella partizione α . Ogni etichetta identifica l'atomo di appartenenza:

$$i \sim j \iff \lambda(i) = \lambda(j)$$

Non si possono fare ipotesi sulla posizione relativa dei siti nè sulle etichette, che formano un vettore di lunghezza L .

2.2.1 Struttura dei vicini

Abbiamo definito il criterio di partizionamento come una relazione di equivalenza *globale*, che non dipende dalla relativa posizione dei singoli siti. In pratica, spesso andiamo a ricostruire la relazione a partire da equivalenze *locali*, definite tra siti “vicini”³, mettendo nello stesso atomo siti “lontani” che sono equivalenti tra di loro tramite una sequenza di equivalenze locali.

Il metodo per partizionare una configurazione generica sullo spazio \mathbf{M} consiste nel raccogliere tutta l'informazione in una matrice dei vicini NNB, concettualmente la versione ridotta (ad un numero opportuno di vicini) e *sparse* della matrice di adiacenza, per poi usare un algoritmo generale di clustering.

In linea di principio, il numero dei vicini per sito non è costante — associamo un vettore di lunghezza variabile sito per sito, a seconda della geometria e configurazione. Per reticoli regolari chiaramente possiamo fissare il numero massimo dei vicini⁴, per ottimizzare il codice, indicando eventualmente con un simbolo convenzionale la fine del vettore (ad es. ci si ferma quando si trova -1). A seconda dell'input, possiamo distinguere vari casi:

- abbiamo in partenza la configurazione completa \mathcal{C} su un reticolo regolare, consideriamo come collegati al sito i -simo tutti i primi vicini j aventi lo stesso simbolo, cioè

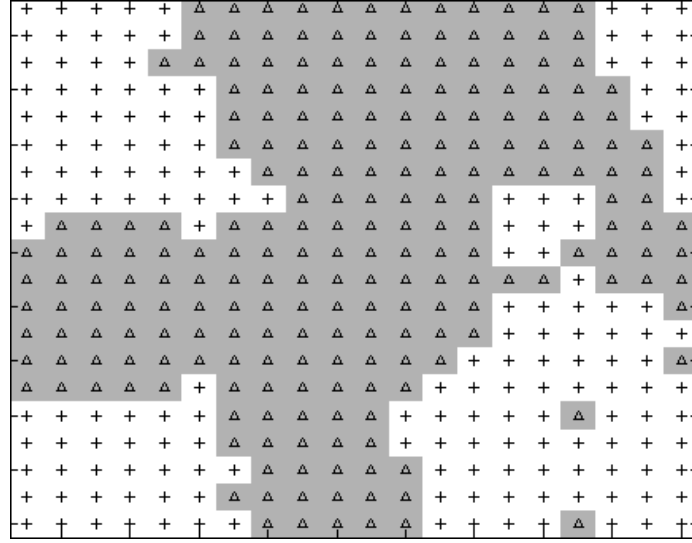
$$d(i, j) = 1, f(i) = f(j)$$

e li aggiungiamo al vettore NNB $[i] [k++] = j$, dove i indica la colonna, k cresce con il numero di j trovati. A questo caso appartiene ad esempio lo studio delle distanze tra reticoli con configurazioni di Ising: si considerano solo i primi vicini aventi lo stesso spin

- nel caso delle partizioni con salto, aggiungiamo a NNB solo il primo sito distante meno di n_s avente lo stesso simbolo, non serve collegare i più lontani (saranno già vicini del precedente)
- non guardiamo affatto le configurazioni, e i vicini sono dati interamente dalla matrice di adiacenza: $\text{NNB}[i] = \text{nonzero}(A[i])$

³Si intende qui nel senso più generale possibile, indipendentemente dalla posizione effettiva, ma legati da una qualunque relazione. Non vi è quindi alcuna perdita di generalità.

⁴Sia per primi vicini, che per un numero allargato, questo è costante.

Figura 2.2.1: Reticolo semplice con $|\mathbb{K}| = 2$

Il caso in cui associamo insieme tutti i simboli dello stesso tipo presenti sul grafo, non è più utile partizionare a partire dall'elenco dei vicini⁵. Si utilizza un diverso algoritmo, una variante del sort descritto a pagina 22, per ottenere una configurazione in cui tutti i simboli uguali sono consecutivi.

Facciamo l'esempio del reticolo bidimensionale ai primi vicini, come si vede in figura 2.2.1. Scorrendo il reticolo in modo ordinato, dall'alto verso il basso, da sinistra verso destra, ovvero secondo la naturale disposizione dei siti in memoria, si vede che per allargare progressivamente (a macchia d'olio) le partizioni, è conveniente considerare i siti vicini già visitati. Per il sito i -esimo si guarda quindi il vicino in alto (*nearest neighbour up*, nnu) e quello a sinistra (*nearest neighbour left*, nnl). In figura 2.2.2 sono disegnate le linee che collegano ciascun sito con il corrispettivo nnl e nnu.

In generale reticoli regolari di dimensione d sono generati da d vettori ad entrate intere: per controllare tutte le relazioni di vicinanza di ogni sito del reticolo è quindi sufficiente esaminare d vicini per sito, invece del numero di coordinazione $z = 2d$. Avere un collegamento tra due siti è una relazione simmetrica, per cui non è necessario controllare in entrambi i versi. Se invece considerassimo solo un sottoinsieme dello spazio \mathbf{M} per avere tutti i collegamenti bisogna effettivamente analizzare z primi vicini.

Problema. Qual è la struttura ottimale (minima) di vicini per riottenere la partizione data?

Basta unire uno-per-uno i siti all'interno di ogni atomo. Ovvero ogni sito ha un solo vicino, il sito che lo precede facente parte dello stesso atomo (o se stesso, se è il primo o unico). Questo è chiaramente sufficiente per ottenere di nuovo la partizione, si può scrivere come un vettore lungo L , indipendentemente dalla geometria del sistema e dal numero effettivo di vicini: è un modo estramamente comodo ed efficiente di catturare le informazioni a partire da NNB, necessarie per la riduzione, ma che possono occupare molto spazio in memoria. Denominiamo l'array contenente il vicino precedente come `previous[i]`.

2.2.2 Partizionamento

Una volta completato l'elenco essenziale dei vicini, si procede al partizionamento. L'algoritmo è noto in letteratura come algoritmo di Hoshen-Kopelman[4], di cui utilizziamo la versione ottimizzata da Mark Newman[5]. Il problema trattato di *clustering* è equivalente a molti altri: bond-percolation, colorazione di grafo generico, ecc.

⁵I vicini sarebbero definiti dalla partizione una volta individuata.

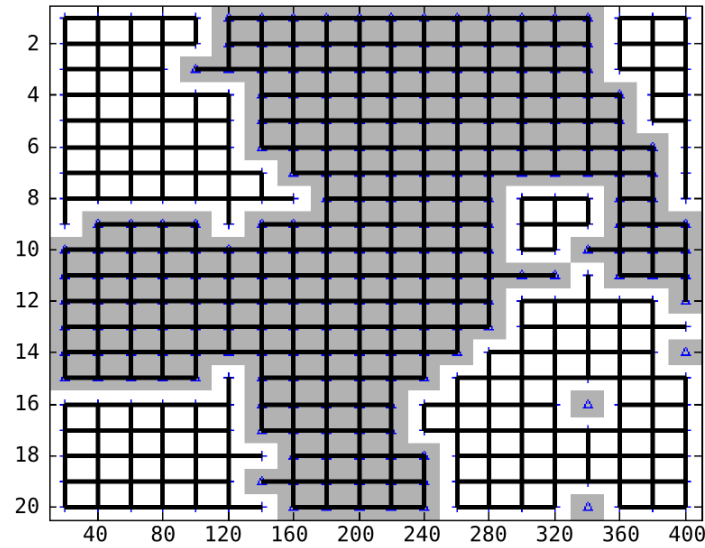


Figura 2.2.2: Collegamento dei siti con stesso simbolo, solo tramite vicino up e vicino left, essendo sufficienti a completare tutti i collegamenti (condizioni al contorno nonperiodiche).

Per il procedimento si crea un array `labels[]`, lungo L , che alla fine dell'algoritmo è codificato nel seguente modo:

`labels[i] < 0` i è la *root* o *radice* del suo cluster, e `-labels[i]` rappresenta la dimensione del cluster – ovvero quanti siti vi appartengono in totale

`labels[i] > 0` il sito punta ad un altro cluster, con *root* alla posizione `labels[i]`

L'array è un insieme di radici e di puntatori ad altri siti – rappresenta cioè un albero, con collegamenti da ogni foglia (sito), alla radice corrispondente, che contiene il numero di siti che le appartengono. Per sapere a che cluster appartiene un qualunque sito, definiamo la funzione `findroot`:

```
int findroot(int i, int labels[]) {
    if (labels[i] < 0) return i;
    return labels[i] = findroot(labels[i], labels);
}
```

definita ricorsivamente sull'albero. Se il sito ha un label negativo, allora è radice. Altrimenti, risaliamo al sito a cui punta e reitro il procedimento. Ad ogni passo si esegue la *path compression*: ogni sito su cui si passa è aggiornato per puntare direttamente alla radice – in questo modo la volta successiva il percorso non deve essere ricostruito.

L'algoritmo è il seguente:

1. ad ogni sito si associa l'atomo banale, formato dal sito stesso, `labels[i] = -1`. In questo modo da subito ad ogni sito corrisponde un cluster, che verrà eventualmente allargato unendolo con i suoi vicini.
2. per ogni sito i , in ordine crescente, controlliamo (fino ad esaurimento) il vettore `NNB[i]`.
 - (a) cerchiamo il cluster di appartenenza, `ri=findroot(i, labels)`
 - (b) per ogni vicino j , identifichiamo a sua volta il cluster, `rj=findroot(j, labels)`
 - (c) se la radice è la stessa, ovvero appartengono già allo stesso cluster, non si fa nulla

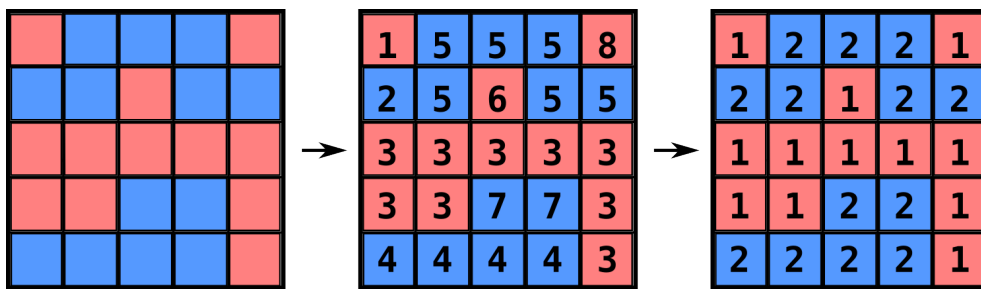


Figura 2.2.3: Funzionamento dell'algoritmo di Hoshen-Kopelman: dalla configurazione indicata dai due colori, si formano i cluster guardando solo ai primi vicini, assegnando una nuova etichetta quando non vi è un primo vicino dello stesso tipo; successivamente si riconoscono globalmente le etichette equivalenti, anche attraverso le condizioni al bordo periodiche.

- (d) se le radici sono diverse, si uniscono i cluster. C'è sicuramente un cluster più grosso (o uguale) e uno più piccolo, rispettivamente di radice r_{\max} e r_{\min} . Per fare il numero minimo di cambiamenti, si aggiunge alla radice il numero di siti accorpati

```
labels[rmax] -= labels[rmin]
```

e si trasforma r_{\min} in una foglia del cluster più grosso:

```
labels[rmin] = rmax
```

La volta successiva che si cerca la root dei siti che appartenevano al cluster più piccolo, si risalirà fino alla nuova radice — il cluster è stato effettivamente assimilato, semplicemente cambiando il valore di un puntatore.

3. alla fine tutti i siti sono stati controllati, $r = \text{findroot}(i, \text{labels})$ punta alla radice di appartenenza, lo spazio è partizionato, l'intero r identifica in maniera univoca le classi di equivalenza.

Riassunto, l'algoritmo corrisponde sostanzialmente a due passaggi, che nel nostro caso avvengono contemporaneamente ma un sito alla volta, illustrati in figura 2.2.3:

- Collegamento dei vicini (FIND)
- Raggruppamento graduale dei cluster così trovati (UNION)

L'algoritmo è ottimale, con runtime $\mathcal{O}(L \ln L)$ e senza bisogno di memoria aggiuntiva. Il fattore logaritmico è dovuto alla fase di unione dei cluster, che in generale richiede $\mathcal{O}(\ln L)$ passi per sito per trovare la radice di appartenenza.

2.2.3 Relabeling

Una volta che abbiamo il vettore `labels[]`, si può utilizzare le informazioni ivi contenute. Si costruisce un vettore di atomi, definendo per ciascuno una struttura che contiene informazioni quali le dimensioni, dove inizia, dove finisce, il *hash* corrispondente.

Per avere il numero degli atomi e inizializzare il vettore corrispondente basta contare quanti siti hanno `labels[i] < 0`. Contemporaneamente possiamo associare ad ogni atomo le sue dimensioni, `-labels[i]`, il suo sito di inizio `i`, e l'ultimo sito trovato, temporaneamente `i`. Ogni elemento del vettore `labels` punta ad un sito (o indica lo stesso come radice) per cui è contenuto nell'intervallo $[0, \dots, L]$ ed è possibile tenere conto di ogni atomo che si incontra scorrendo tutti i siti e utilizzare questo contatore come nuova etichetta: non più la radice, ma il numero dell'atomo corrispondente in $[0, \dots, n]$.

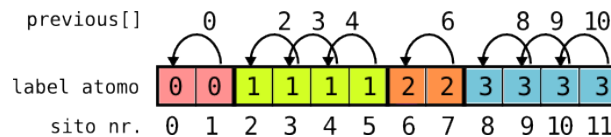


Figura 2.2.4: Per ogni sito è stato indicato l'atomo di appartenenza, sotto l'indice, sopra il valore del vettore che punta al sito precedente appartenente allo stesso atomo.



Figura 2.2.5: Sintetica rappresentazione della struttura in memoria che rappresenta le quantità collegate ad ogni sito

Avendo un elenco di atomi contenente le dimensioni di ciascuno, si calcola l'entropia della partizione in modo diretto, sommando le dimensioni dei vari atomi moltiplicate per il loro logaritmo e normalizzando.

Generiamo inoltre il vettore ottimale dei vicini: `previous[]`. Per ogni sito i , conosciamo infatti l'atomo di appartenenza e il precedente sito trovato appartenente allo stesso atomo: `previous[i]` conterrà proprio questo dato. Aggiorniamo poi il puntatore all'ultimo sito trovato con i stesso. Andando a leggere il vettore `previous[]` in modo ordinato, abbiamo una lista ordinata, che ad ogni intero associa un numero strettamente minore o uguale (se il sito non ha un precedente).

2.2.4 Hashing

Definiamo un *hash*⁶ come un modo univoco di associare ad ogni atomo un elemento dello spazio delle stringhe a 64bit — il numero degli atomi è sicuramente molto inferiore a 2^{64} per cui con un buon algoritmo è improbabile che due atomi diversi ottengano lo stesso identificativo[6][7].

La conoscenza di un atomo corrisponde alla lista di siti che lo compongono, tuttavia sarebbe oneroso tenere da parte un'altra volta questa lista, ordinata stavolta per atomo, oppure andarla a ricostruire ogni volta. Per controllare l'uguaglianza tra due atomi di due diverse partizioni, bisognerebbe fare un confronto sito per sito. Se ogni atomo diverso è distinguibile solo con il confronto di un intero (operazione atomica), il confronto tra due atomi è immediato e $\mathcal{O}(1)$ in tempo, indipendentemente dalle dimensioni dell'atomo.

La chiave sta nel trovare un buon algoritmo di hashing, che sparpagli i bit di informazione contenuti dall'atomo in maniera uniforme nello spazio delle stringhe binarie a 64bit, in modo che a input diverso corrisponda *sempre* (a meno di una probabilità idealmente prossima a 2^{-64}) un identificativo diverso. L'algoritmo utilizzato in questo caso è l'estremamente semplice *hash DJB*⁷, dopo un confronto con alcuni altri algoritmi[8], risultati dare molte collisioni, ovvero capitava che triple di interi diverse dessero risultati uguali[9].

Gli atomi sono equivalenti ad array di interi, per cui l'algoritmo prende un array di lunghezza arbitraria di interi a 32bit e restituisce un intero a 64bit, partendo da un valore costante iniziale, fissato empiricamente a 5381. L'aggiunta di un intero `value` alla variabile `hash` avviene nel modo seguente:

⁶Etimologicamente, *to hash* vuol dire maciullare, sparpagliare

⁷Acronimo dell'ideatore, D. J. Bernstein

```

void DJBHash_step(u_int64_t &hash, u_int32_t value) {
    hash = (hash << 5) + hash + (0xff000000 & value);
    hash = (hash << 5) + hash + (0x00ff0000 & value);
    hash = (hash << 5) + hash + (0x0000ff00 & value);
    hash = (hash << 5) + hash + (0x000000ff & value);
}

```

ovvero la variabile a 32bit che codifica il valore del sito viene spezzata in 4 pezzi da 1 byte ciascuno; ogni byte viene sommato al hash, mentre questo viene rimischiato tramite uno shift e una somma del valore precedente. Senza spezzettare i 4 byte il rimescolamento dei bit risulta insufficiente e il hash perde di unicità.

Il hash viene calcolato intanto che si scorre l'array dei labels e si costruisce il vettore `previous`. Per ogni sito i recuperiamo l'indice dell'atomo corrispondente `labels[i]`, che contiene una variabile preposta all'immagazzinamento del hash, al quale aggiungiamo il valore di i di volta in volta.

2.3 Operazioni tra partizioni generiche

2.3.1 Prodotto

Dovendo ragionare sui label assegnati ai siti, cerchiamo un modo di labellare in modo univoco i siti nella partizione prodotto, con un'operazione sito per sito, ovvero locale. Dalla definizione si ottiene che un sito i appartiene all'atomo $(\alpha \vee \beta)_{jk}$ se appartiene a $A_j \cap B_k$, ovvero se ha $\lambda_\alpha(i) = j$ e $\lambda_\beta(i) = k$.

Il prodotto deve mantenere la topologia, quindi prendendo due siti "vicini", ovvero appartenenti allo stesso atomo in entrambe le partizioni, devono rimanere nello stesso atomo del prodotto

$$m, n \in A_j \text{ e } m, n \in B_k \implies m, n \in (\alpha \vee \beta)_{jk}$$

ma i siti m, n hanno label $\lambda_\alpha(m) = \lambda_\alpha(n) = j$ e $\lambda_\beta(m) = \lambda_\beta(n) = k$ rispettivamente per le due partizioni — devono avere lo stesso label nel prodotto, che può dipendere solo dal label degli atomi di partenza. L'unica scelta possibile è assegnare label (j, k) (o equivalente) ad ogni sito del prodotto, ovvero la coppia (ordinata) dei label di partenza. Una volta eseguito l'assegnamento, gli atomi sono riconosciuti cercando label dello stesso valore. Tutto il procedimento è illustrato in figura 2.3.1.

Se i label sono rappresentati da interi a 32bit, la coppia può essere rappresentata da un intero a 64bit, avente nei primi 32bit il valore del primo label, e negli ultimi quello del secondo: in questo modo tutte le operazioni sono atomiche su un moderno processore.

Ricapitolando: al sito i -esimo, partendo da $\lambda_\alpha(i) = j$ e $\lambda_\beta(i) = k$, si assegna nel prodotto $\lambda_{\alpha\vee\beta}(i) = (j, k)$:

```
prod[i] = alpha[i] << 32 | beta[i];
```

come si vede dalla formula esplicita, $(j, k) \neq (k, j)$ in quanto gli indici identificano atomi appartenenti a partizioni diverse, con siti in genere diversi — se al posto del OR binario si volesse usare un'altra operazione, va cercata una asimmetria nei due indici, per non unire atomi distinti. Ad esempio una partizione con labels $\{1, 2\}$ è equivalente a $\{2, 1\}$ (sempre due atomi distinti); il prodotto è ancora una partizione con due atomi, ma le etichette del prodotto sono: $\{(1, 2), (2, 1)\}$ — devono essere chiaramente diverse.

Anche in questo caso il prodotto comporta un numero di operazioni sito-per-sito, con complessità $\mathcal{O}(L)$.

2.3.2 Entropia

Per calcolare il numero di siti di un atomo, ovvero corrispondente ad ogni etichetta univoca, bisogna innanzitutto scorrere l'elenco per cercare tutte le etichette e ogni volta che

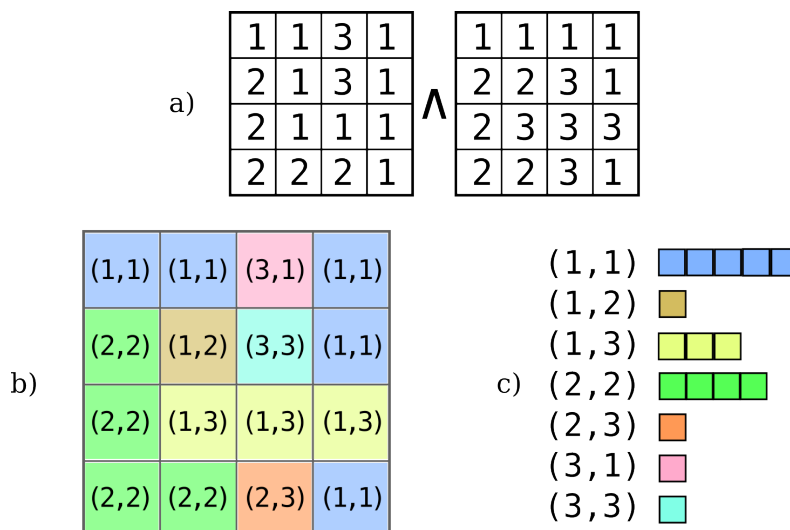


Figura 2.3.1: Calcolo di un prodotto tra partizioni. a) partizioni di partenza; b) partizione prodotto; c) lunghezze degli atomi ordinati del prodotto, al fine di calcolare l'entropia

si trova una non precedentemente catalogata, scorrere l'intero array in avanti e contare quante volte appare. Questa è chiaramente un'operazione con complessità

$$\sum_{i=1}^L (L-i) = \mathcal{O}(L^2)$$

Ma si può fare meglio. Infatti non è necessario scorrere N volte l'intero elenco dei labels. Si può riordinare questo elenco tramite un'operazione di *sort*⁸. L'andamento asintotico è noto in letteratura essere $\mathcal{O}(L \log_2 L)$ nel numero di siti, con una costante di proporzionalità dipendente dall'algoritmo, ma non particolarmente onerosa.

Ordinando l'elenco delle etichette (non importa neanche se in ordine crescente o decrescente) si ottiene un array in cui sicuramente tutti gli elementi uguali sono consecutivi – si procede a questo punto al calcolo dell'entropia come se si trattasse di una sequenza unidimensionale.

L'altro metodo, computazionalmente più semplice, $\mathcal{O}(L)$ sia in runtime che in memoria, che consiste nel numerare temporaneamente tutti i label trovati in ordine crescente e rietichettarli con la successione $\{1, 2, \dots\}$ non è gestibile, in quanto i label possono essere distribuiti in modo “sparso” nello spazio delle coppie (j, k) . In generale vi possono essere al massimo L atomi nel prodotto, ma $L \times L$ coppie, quindi sicuramente non tutte sono realizzate. La complessità dell'organizzarle in qualche altro modo, ad esempio in un albero binario bilanciato o un hash – tutte superano in fretta la semplicità del *sort* e anche nel caso migliore hanno lo stesso andamento asintotico. In pratica l'andamento è polinomialmente peggiore, per motivi di mal gestione della cache L1 e L2 del processore, in quanto sia gli alberi che i hash risultano nonlocali, mentre un *sort* sposta gli elementi in modo ordinato e accede alla memoria in modo più lineare.

Una volta calcolata l'entropia di ogni partizione, sia di Shannon che topologica, la distanza di Rohlin segue immediatamente dalla formula 1.1.3 a pagina 7, senza nessuna modifica.

2.3.3 Differenza simmetrica

Per calcolare la differenza tra due atomi, illustriamo come calcolare in modo efficiente la differenza simmetrica tra due insiemi.

⁸Nello specifico ci affidiamo al *quicksort* delle librerie glibc

A1 = {9, 8, 5, 4, 1, 0}
A2 = {8, 7, 5, 4, 2, 1}
allineamento1 9 8 - 5 4 3 - 1 0
allineamento2 - 8 7 5 4 - 2 1 -
differenza? + . + . . + + . + == 5

Figura 2.3.2: Differenza simmetrica tra insiemi ordinati, rappresentata testualmente

Dati due insiemi ordinati di dimensione n e m , è possibile calcolare la differenza simmetrica in modo lineare, con al massimo $2(n + m + 1)$ operazioni, mentre per due insiemi disordinati sono necessarie $\mathcal{O}(nm)$ operazioni. Il procedimento è illustrato in figura 2.3.2.

La prima fase è quella di allineamento, ovvero si scorrono i due array rappresentanti gli insiemi, aggiungendo spazi all'uno o all'altro, in maniera tale che i valori nonvuoti siano uguali nelle due righe. Successivamente si contano le differenze tra le due stringhe di simboli allineati.

Chiaramente si possono fare entrambe le operazioni contemporaneamente, contando +1 ogni volta che si aggiunge uno spazio. Quando si vuol sapere se due atomi hanno al massimo n_ϵ siti diversi, ci si può fermare dopo aver inserito n_ϵ spazi. Non è poi neanche necessario rappresentare gli spazi in memoria, basta scorrere in avanti il puntatore corrispondente per allinare al volo intanto che si leggono gli insiemi stessi.

L'operazione è estremamente efficiente se si ha l'elenco corrispondente ai siti dei due insiemi. Nel nostro caso, a partire da un qualunque sito, è sufficiente usare il vettore `previous[]` per ottenere il sito precedente, sicuramente ordinato.

2.3.4 Intersezione

Definiamo ora il modo operativo di calcolare la partizione comune, $\sigma = \alpha \wedge \beta$. Descriviamo il processo prima in termini di atomi, poi in termini di siti, per arrivare all'algoritmo ottimale.

La partizione intersezione, σ essendo la partizione più fine tale che $\sigma \leq \alpha$ e $\sigma \leq \beta$, ha atomi formati dall'unione minimale di atomi sia di α che di β . Ovvero ogni atomo $C_k = \cup_j A_j = \cup_m B_m$. Chiaramente andare a cercare le combinazioni di A_j e B_m necessarie è un'impresa disperata.

In termini di siti, vediamo che ogni atomo C_k collega siti appartenenti a diversi atomi di α se questi sono collegati in β . Questo corrisponde già ad un utile criterio:

$$i \sim j \text{ in } \sigma \iff i \sim j \text{ in } \alpha \vee i \sim j \text{ in } \beta \quad (2.3.1)$$

che diventa ancora più utile se ci si ricorda che la relazione di equivalenza tra siti è mantenuta a livello locale dalla struttura NNB precedentemente introdotta. Ovvero possiamo esprimere la relazione di equivalenza 2.3.1 a livello locale, sito per sito, ottenendo la corretta partizione globale alla fine!

La partizione σ ha anch'essa una struttura di vicini corrispondenti:

- i è vicino di j anche se è solo vicino in α
- i è vicino di j anche se è se è vicino in β

la partizione σ sarà quindi ricostruita tramite un'opportuna struttura NNB, che ha in generale il doppio dei vicini $\tilde{z}(\sigma) = \tilde{z}(\alpha) + \tilde{z}(\beta) = 2\tilde{z}$, ovvero sito-per-sito l'unione dei vicini di α e β .

- $\text{NNB}_\sigma[i][2k] = \text{NNB}_\alpha[i][k]$
- $\text{NNB}_\sigma[i][2k+1] = \text{NNB}_\beta[i][k]$

dove per NNB_α si intende la struttura corrispondente utilizzata per il partizionamento di α , $0 \leq k < \bar{z}$.

Una volta definita la struttura dei vicini, si partiziona l'array contenente tutti i siti con l'algoritmo descritto a pagina 18.

Ottimizzazione Non è in realtà necessario aggiungere i vicini di ogni sito, come definiti originariamente dalla struttura NNB corrispondente. Esiste un vettore di vicini ottimale, che ricostruisce esattamente le partizioni originarie: `previous []`.

Il problema dell'intersezione è quindi ridotto a considerare 2 vicini per ogni sito, invece che $2d$ senza quest'ottimizzazione! Per un qualunque sito definiamo quindi

- $NNB_\sigma[i][0] = \text{previous}_\alpha[i]$
- $NNB_\sigma[i][1] = \text{previous}_\beta[i]$

2.4 Riduzione

Una volta che si è in grado di calcolare l'intersezione, è possibile calcolare le partizioni ridotte. Vediamo ora come implementare i diversi criteri esposti a pagina 8. Lo scopo principale è quello di evitare la costruzione, il prodotto e l'intersezione dei fattori dicotomici, per motivi di tempo e occupazione in memoria.

Visto che le partizioni ridotte servono solo al calcolo del loro prodotto, costruiremo per ciascuno i giusti labels. Non serve costruire i vettori `previous []`, nè le strutture degli atomi, ma solo calcolare l'entropia a partire dai labels con il metodo presentato a pagina a pagina 22 visto che è necessaria per la distanza di Rohlin.

2.4.1 Confronto con l'intersezione

Per ridurre due partizioni α e β definiamo innanzitutto l'intersezione $\sigma = \alpha \wedge \beta$ calcolata con l'algoritmo appena presentato. In generale per fare molti confronti, sarà necessario usare una variabile temporanea che contenga l'intersezione di volta in volta.

Costruiamo le partizioni dicotomiche semplici, ovvero a partire da ogni atomo $A_k \in \alpha$, definiamo $\hat{\alpha}_k = \{A_k, A_k^c\}$:

Costruzione in memoria si tiene l'elenco degli atomi, ovvero un array di strutture `atomi []`, su cui possiamo scorrere. Per ogni atomo (ad esempio il k -esimo) conosco il sito finale, `atomi[k].end`, ricostruisco gli altri siti appartenenti tramite iterazione del vettore `previous[atomi[k].end]`. Ho inoltre il sito di inizio, `atomi[k].start`.

Confrontiamo ogni partizione dicotomica $\hat{\alpha}_k$ con l'intersezione, $\rho = \hat{\alpha}_k \wedge \sigma$. Se $\rho \neq \nu$, la partizione banale, scarto la partizione dicotomica $\hat{\alpha}_k$. Si hanno i possibili casi:

Intersezione nonbanale $\exists C_m \in \sigma$ tale che $A_k = C_m$: in questo caso l'atomo di α corrisponde esattamente ad un atomo di σ , ovvero hanno tutti i siti in comune. La loro intersezione conterrà almeno l'atomo $\hat{\alpha}_k$, per cui sarà non banale \implies scarto A_k .

Intersezione banale Se A_k non corrisponde esattamente a un qualche C_m , allora σ unisce più atomi di α . In questo caso, la partizione σ collega siti di A_k con siti del suo complementare, banalizzando l'intersezione con la partizione dicotomica $\implies \hat{\alpha}_k \wedge \sigma = \nu$.

Come si calcola se esiste esattamente un atomo C_m corrispondente? Visto che stiamo iterando sull'elenco degli atomi di α , conosciamo sicuramente un sito appartenente ad A_k che chiamiamo j , `j=atomi[k].end` e l'atomo corrispondente al sito j nella partizione σ , tramite il vettore delle etichette di σ . Fatto questo, abbiamo le strutture degli atomi corrispondenti al sito j per entrambe le partizioni, per cui conosciamo il *hash* corrispondente a C_m e a A_k . Se il *hash* è uguale, allora gli atomi sono uguali! Per il confronto del fattore dicotomico generato dall'atomo A_k abbiamo quindi in totale fatto 4 operazioni veloci:

1. Dall'indice k si trova la struttura contenente un sito j e il hash, $\text{atoms}[k] \rightarrow 1$ lookup
2. Si cerca l'indice dell'atomo corrispondente in $\sigma, \text{labels}_\sigma[j] \rightarrow 1$ lookup
3. Si trova la struttura corrispondente all'atomo in σ e si legge il *hash* $\rightarrow 1$ lookup
4. Si confrontano i due hash $\rightarrow 1$ confronto

Rimane da moltiplicare insieme tutte le partizioni dicotomiche rimaste, per ottenere la partizione ridotta $\hat{\alpha} = \bigvee_k \hat{\alpha}_k$. Numeriamo i fattori dicotomici rimasti con un indice crescente, a partire da 2. Vediamo come in ogni partizione dicotomica vi sia un atomo della partizione originaria e il suo complementare: gli atomi sono disgiunti, per cui il prodotto dei singoli atomi risulta vuoto. Rimangono da moltiplicare tutti i complementari per gli atomi rimasti: genero nuovamente esattamente gli atomi che non ho scartato nella riduzione e un ulteriore atomo (generalmente sconnesso) composto dall'unione di tutti gli atomi scartati.

Ricostruzione Creiamo una partizione banale con tutti i label pari a 1, a rappresentare l'atomo di "background". Mettiamo un label $k \geq 2$ ai siti di ogni atomo A_k che non è stato scartato nel confronto. Questo è semplice, basta porre uguali a k tutti i siti partendo da $j = \text{atoms}[k].\text{end}$ e iterando tramite $j = \text{previous}[j]$ fino a che non arriviamo all'inizio dell'atomo. Alla fine abbiamo una partizione in cui l'atomo disconnesso con gli scarti ha label 1, mentre gli altri siti hanno label numerati a partire da 2. Il processo è chiaramente proporzionale al numero di siti che abbiamo dovuto riscrivere, quindi

$$\mathcal{O}(L) + \mathcal{O}\left(\sum_k \mu(A_k)\right) = \mathcal{O}(L)$$

Abbiamo costruito in tempo lineare la partizione ridotta, in quanto per ogni atomo c'è un confronto atomico e il relabeling del risultato corrispondente.

2.4.2 Confronto diretto

In questo tipo di riduzione, confrontiamo direttamente gli atomi di due partizioni, α e β , scartando le partizioni dicotomiche corrispondenti, nel caso in cui gli atomi soddisfino il criterio prescelto.

Per evitare problemi di definizione, consideriamo la partizione β fissata, costruendo $\hat{\alpha}$ a partire dal confronto di α e β .⁹

Abbiamo due criteri di confronto:

1. Scartiamo l'atomo A_k se è uguale a $B_m \in \beta$ definito sugli stessi siti (usando il confronto tra hash)
2. Scartiamo l'atomo A_k se differisce da $B_m \in \beta$ corrispondente meno di ϵ in misura:

$$\mu(A_k \triangle B_m) \leq \epsilon$$

Si procede molto similmente al caso del confronto con σ , in quanto iteriamo su tutti gli atomi di α , troviamo i siti corrispondenti all'atomo cercato, da questi gli atomi di β sopra quei siti e eseguiamo il confronto.

La differenza simmetrica è stata definita in termini di numero di siti diversi n_ϵ , mentre sopra in termini della misura normalizzata che deve essere $\leq \epsilon$, definiamo dunque:

$$\epsilon = \frac{n_\epsilon}{L}$$

La partizione $\hat{\alpha}$ è ricostruita a partire dagli atomi rimasti come spiegato nella sezione precedente. L'unica differenza è insomma che il confronto viene fatto con gli atomi corrispondenti di β e non di σ e che il criterio non è l'uguaglianza, ma l'uguaglianza a meno di ϵ (che può essere 0).

⁹Poi si fa il contrario, tenendo fisso α , generando $\hat{\beta}$.

Capitolo 3

Proteine dei virus dell'influenza

Esponiamo in questo capitolo i risultati ottenuti dalle distanze tra sequenze di proteine, ottenute dai database biologici, dei virus dell'influenza A umana, di tipi H1N1 e H3N2.

Lo studio delle sequenze comprende fasi separate:

1. Raccolta delle sequenze e preanalisi
2. Compilazione matrice delle distanze tra le sequenze
3. Clustering a partire dalle distanze
4. Estrazione informazioni dai cluster e predizioni

Introduciamo innanzitutto il virus dell'influenza e alcune sue caratteristiche fondamentali, che motivano la scelta di che cosa analizzare con i metodi entropici. Successivamente esaminiamo nel dettaglio i passi elencati sopra.

3.1 Il virus dell'influenza

I virus dell'influenza sono caratterizzati da genomi segmentati in forma di RNA a singolo filamento negativo, dipendenti da una RNA-polimerasi di origine virale per la trascrizione prima della replicazione[10].

I virus sono dei tre tipi A, B e C, rappresentanti tre dei cinque generi della famiglia *Orthomyxoviridae*, caratterizzati da genomi segmentati. I virus hanno un'origine genetica comune, tuttavia la divergenza evolutiva ha reso impossibile il riassortimento di materiale genetico tra i tre diversi tipi.

I virus dell'influenza B sono caratterizzati da un minore tasso di mutazione (circa $1/3$ rispetto all'influenza A) e attaccano quasi esclusivamente l'uomo, quindi non tendono a scatenare pandemie. I virus del tipo C hanno un tasso di mutazione ancor più lento e sono asintomatici, per cui concentriamo la nostra attenzione sul virus A, il più pericoloso e fonte di annuali epidemie – definite come trasmissioni prolungate, diffuse da persona a persona.

3.1.1 Influenza A

La forma del virus A è sferica, dal diametro di circa 100 nm, ricoperta di numerose proteine, circa 500 molecole di HA e 100 di NA, che protudono dalla superficie. Sono in particolare *antigeni*, termine che vuol dire **generatore** di **anticorpi**, in quanto il sistema immunitario riconosce e reagisce al virus proprio tramite gli antigeni.

All'interno il virus contiene 8 segmenti di RNA, che codificano le proteine per la replicazione. Vediamo ora nel dettaglio le funzioni delle proteine di superficie:

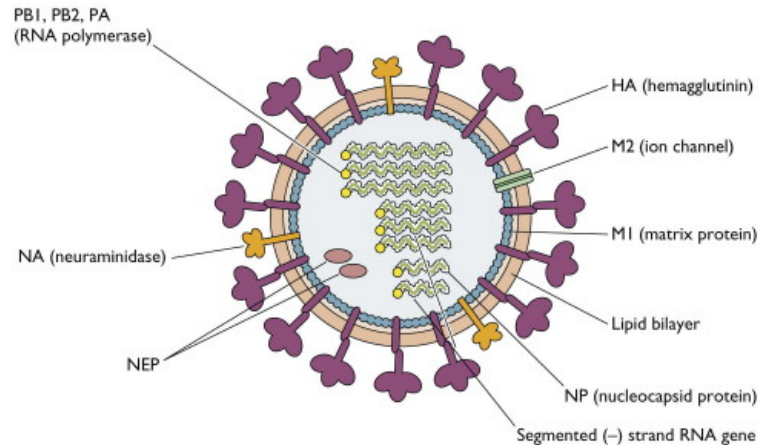


Figura 3.1.1: Forma del virus dell'influenza A, in cui sono evidenziati gli antigeni all'esterno e la disposizione del materiale genetico negli 8 segmenti.

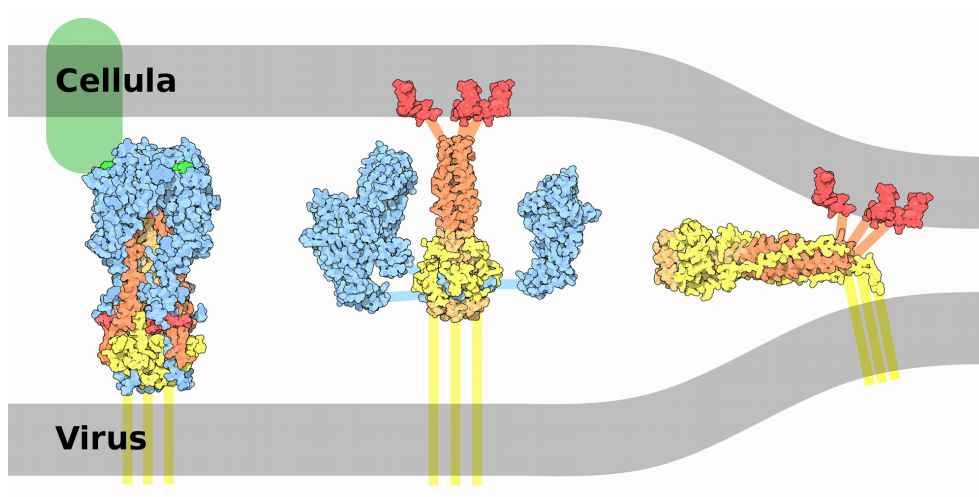


Figura 3.1.2: Illustrazione della molecola HA che protendendo dalla superficie del virus, aggancia la cellula e rende possibile l'iniezione del materiale genetico.

HA Emoagglutinina, proteina responsabile per il binding del virus alla cellula infettata. L'emoagglutinina forma delle "punte" sulla superficie del virus che servono per agganciarsi alla cellula ospite. L'aggancio è necessario per l'efficiente trasferimento del materiale genetico all'interno della cellula, un processo che può essere bloccato dagli anticorpi che si combinano all'emoagglutinina. La forma tridimensionale e il processo con cui la molecola attacca la cellula sono mostrate in figura 3.1.2[11].

NA Neuraminidase, proteina enzimatica che aiuta il rilascio della progenie virale dalle cellule infette.

Le rimanenti proteine M formano il capsido (l'involucro del virus), le proteine NP e NS rilasciano il materiale nella cellula, le proteine PA e PB gestiscono la replicazione del RNA virale.

I virus di tipo A sono classificati in base ai possibili tipi di proteine di superficie: almeno 16 tipi di HA e 9 tipi di NA sono noti al momento della scrittura. Tutte le combinazioni sono espresse negli uccelli, che si considera essere gli originatori dell'influenza e di aver propagato il virus alle altre specie. Le combinazioni sono indicate dalle sigle HxNy, dove x e y sono gli indici corrispondenti, ad esempio:

H1N1 Virus di questo sottotipo hanno causato le pandemie nel 1918 e 2009

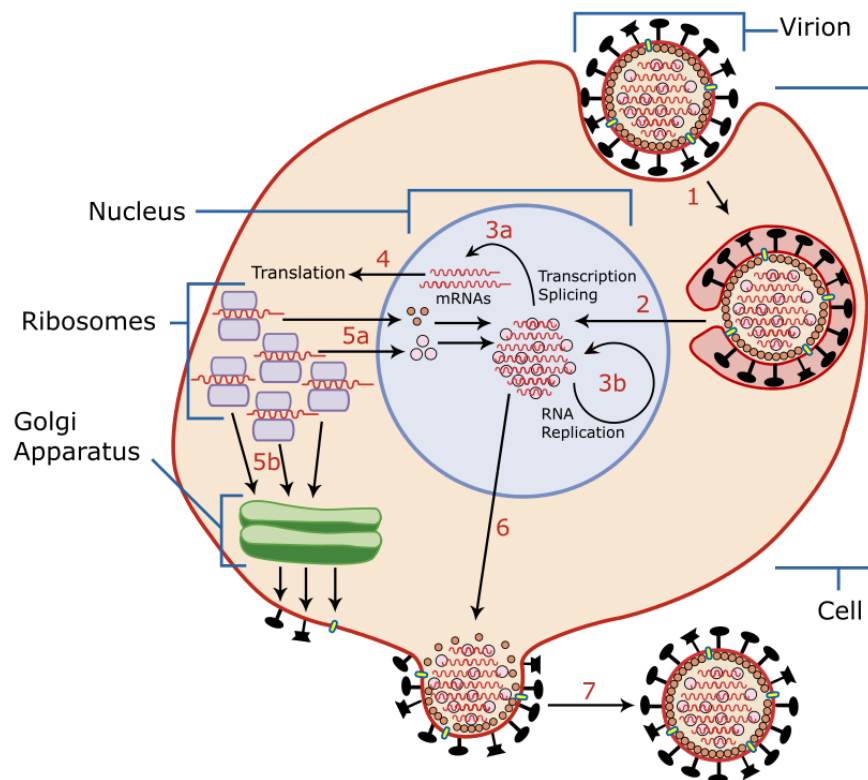


Figura 3.1.3: Processo di riproduzione del virus, in tutte le sue fasi.

H3N2 La “solita” influenza stagionale

H5N1 Una recente variante di origine aviaria dalla forte virulenza, poco diffusa in occidente

Ogni sequenza campionata nei database biologici, è caratterizzata dalla stringa *strain*, ad esempio

A/Human/Fujian/411/2002 (H3N2)

indicante che si tratta di un virus del tipo A, sottotipo H3N2, isolato a Fujian nel 2002 da un essere umano (l'indicazione di provenienza da umano è di solito tralasciata), ed è la 411 sequenza campionata in quell'anno in quel laboratorio.

La pressione evolutiva è concentrata sulle proteine HA e NA, in quanto obiettivi del sistema immunitario. La particolare struttura del genoma del virus dell'influenza e la funzione delle proteine virali consentono due tipi di evoluzione per evadere le risposte immunitarie adattative a lungo termine in vari ospiti:

antigenic drift che corrisponde al normale processo di mutazione puntuale del materiale genetico. La replicazione è affidata all'RNA, che non supporta meccanismi di correzione dell'errore, con in media un errore di trascrizione per ogni replicazione del virus, il che vuol dire che ogni singolo virus è diverso dagli altri. L'enorme tasso di mutazione è la causa della rapida obsolescenza dell'immunità acquisita, che di solito non si estende da una stagione all'altra.

antigenic shift corrispondente alla ricombinazione del materiale genetico provenienti da virus diversi o da sottospecie diverse. Questo accade quando una singola cellula è infettata contemporaneamente da più virus. Il virus dell'influenza ha svariati segmenti indipendenti, che si possono rimescolare all'interno della cellula infetta. I

virus risultanti avranno materiale genetico e proteico fornito dalla combinazione casuale di quelle dei virus che hanno attaccato la cellula. In questo modo da un virus H5N1 e H3N2 si possono formare combinazioni come H3N1 e H5N2 (oltre alle combinazioni di tutte le altre proteine all'interno del virus, non rilevate dal sistema immunitario).

L'*antigenic shift* è la mutazione più pericolosa, in quanto cambia radicalmente l'aspetto del virus per il sistema immunitario, rendendolo irriconoscibile e quindi non più attaccato dagli anticorpi. La pericolosità dello shift è che il virus A dell'influenza non attacca solo umani, ma anche maiali, uccelli, cavalli, pipistrelli e moltissimi altri tipi di animale. L'evoluzione di ciascun tipo procede quindi indipendentemente, fino a che tipi diversi di influenza non attaccano lo stesso individuo e ricombinano i materiali genetici.

Questo è quanto è accaduto nel caso dell'influenza suina. I maiali possono ammalarsi, con sintomi simili all'uomo, di influenze aviarie, suine ed umane contemporaneamente. In particolare la pandemia del 2009, del sottotipo H1N1, è stata causata dalla ricombinazione di sequenze di questi 3 tipi, con antigeni tipici degli uccelli e altre proteine di provenienza mammifera, permettendo al virus di attaccare anche gli esseri umani, il cui sistema immunitario era impreparato al nuovo tipo di virus, anche con trasmissione da uomo a uomo, scatenando una pandemia.

3.1.2 Vaccini

I primi vaccini contro l'influenza sono stati introdotti negli anni 40 del secolo scorso. La vaccinazione ha un effetto statisticamente rilevante nel prevenire un grande numero di casi d'infezione, ridurre i costi e le morti causate dalle complicazioni dell'influenza, presentando un rapporto costo-benefici positivo per tutte le categorie di popolazione definite "a rischio" per le quali viene raccomandata. Il vaccino conferisce però limitata *cross-immunity*, ovvero immunità per ceppi diversi da quello preso ad obiettivo dal vaccino — ciò comporta la necessità di una vaccinazione quasi ad ogni stagione influenzale[12].

Dato la grande variabilità antigenica dei virus dell'influenza, in media la composizione dei vaccini è cambiata ogni 2-5 anni a partire dal 1972. La composizione raccomandata dei vaccini per l'emisfero nord è decisa a febbraio, per l'emisfero sud a settembre, circa mezzo anno prima dell'epidemia stagionale, rispettivamente a settembre e a marzo, in modo da precedere l'epidemia e garantire la massima copertura per la popolazione.

I focolai di epidemie localizzate, della durata media dai 3 ai 6 mesi, generano la maggior parte delle nuove sequenze virali tramite antigenic drift[13], per poi diffonderle nel resto del mondo. L'effetto del drift è difficile da individuare, per il rimescolamento dei virus nella popolazione, per via degli shift, per bias nel processo di campionamento, ecc. Non sia se sia la somma graduale di mutazioni nel tempo accumulate durante le epidemie annuali, o sia un processo puntuato e stocastico che occorre al presentarsi di una mutazione che sfugge al controllo dell'immunità acquisita. La popolazione dei virus è sottoposta ad un collo di bottiglia evolutivo nel periodo estivo, durante il quale la popolazione di organismi infetti è minima e determina il ceppo principale della stagione successiva. È di notevole interesse scientifico e medico individuare quindi al più presto i ceppi dominanti per la stagione epidemica.

Le raccomandazioni per il vaccino sono compilate dal WHO(author?) [14], tramite il suo Global Influenza Surveillance Network. La rete comprende 4 WHO Collaborating Centres (WHO CCs) e 137 istituti in 103 stati, riconosciuti dal WHO come Centri Nazionali per l'Influenza (NIC). I NIC raccolgono campioni nel loro paese, eseguono l'isolamento del virus e una caratterizzazione antigenica preliminare. Inviano poi i nuovi ceppi isolati per l'analisi genetica e antigenica di alto livello alle sedi del WHO, i cui risultati formano la base delle raccomandazioni.

Il vaccino è trivalente, protegge contro 3 diversi ceppi virali. Per tutte le stagioni il vaccino comprende gli stessi tipi, ma diversi ceppi. Per la stagione 2012-2013 le raccomandazioni sono le seguenti(author?) [14]:

1. un sottotipo A/H1N1, A/California/7/2009

2. un sottotipo A/H3N2, A/Victoria/361/2011
3. un tipo B, B/Wisconsin/1/2010

Il ceppo è scelto tra i virus dello stesso sottotipo, da un'analisi antigenica *in vitro*. I campioni provenienti dai NIC vengono classificati in base alla potenza della reazione con l'antisiero sviluppato a partire dai ceppi candidati al vaccino della stagione successiva. Il ceppo che ha un numero maggiore di riscontri viene incluso nel vaccino. Il criterio è sempre in ritardo, in quanto non è in grado di predire la prevalenza dei virus nella stagione successiva, ma solamente indicare qual è la più diffusa nella stagione precedente la raccomandazione. Negli ultimi anni questo ha causato 6 errori su 19, dove per errori si intende la non inclusione del ceppo dominante la stagione epidemica. In particolare la raccomandazione per la stagione 2011-2012 era:

1. un sottotipo A/H1N1, A/California/7/2009
2. un sottotipo A/H3N2, A/Perth/16/2009
3. un tipo B, B/Brisbane/60/2008

la raccomandazione per A/Perth è risultata inappropriata (pur con l'effetto mitigante di una parziale cross-immunity), in quanto un numero maggiore di casi è stato riportato durante la stagione epidemica per il ceppo A/Victoria, come predetto dall'algoritmo in [15].

3.2 Clustering

Il *clustering* è definito come il modo di raggruppare oggetti nello stesso cluster sono più simili tra loro che a oggetti appartenenti ad altri cluster. Gli utilizzi sono molteplici: riconoscere oggetti come simili nonostante le differenze, distinguere il primo piano dallo sfondo, ottenere una efficiente compressione utilizzando solo l'indicativo del cluster di appartenenza, ecc. In particolare ci poniamo lo scopo di associare ogni sequenza virale ad un cluster, in maniera da poter parlare di k ceppi invece di N sequenze, con $k \ll N$, che è il modo più naturale di ragionare sul problema.

Una clusterizzazione è una partizione di un insieme, in cui ad ogni elemento i viene associato in modo rigido a un cluster C di appartenenza:

$$i \rightarrow C_i$$

attraverso una opportuna scelta di un criterio, basandosi sui dati ottenuti, sulle informazioni aggiuntive che si possono fornire, sul tipo di elementi dell'insieme. In generale il problema del clustering, ovvero la disposizione degli elementi in k cluster in modo tale da minimizzare una funzione globale di costo, è NP completo per la natura combinatoria del problema. Vi sono molti metodi euristici per ottenere soluzioni approssimate, di cui noi abbiamo verificato l'applicabilità, ma una descrizione completa è argomento di ricerca a sé stante [16].

L'algoritmo più semplice e utilizzato di questi è *kmeans*, che deve il nome al numero k di cluster da generare come parametro di input. Si richiede che i punti siano embedded in uno spazio vettoriale, rendendo difficile l'applicazione, dato che abbiamo solo una matrice di distanze. L'algoritmo cerca di separare gli elementi in gruppi di uguale varianza, minimizzando il "momento d'inerzia" definito come la somma delle distanze dall'elemento centrale di ogni cluster. Per ogni cluster si raggruppano gli elementi attorno a dei centri scelti inizialmente a random (di qui l'assunzione dell'embedding in \mathbb{R}^n), per poi spostare i punti da un cluster all'altro e ricalcolando i centri, fino al momento d'inerzia totale non si stabilizza attorno ad un minimo. Tuttavia partendo da una soluzione random si raggiunge ogni volta un minimo e una clusterizzazione finale diversa, spesso con pessimi risultati. Un modo innovativo di generare l'embedding e garantire un buon comportamento di *kmeans* è presentato nella sezione sui metodi spettrali.

Introduciamo ora il metodo principale utilizzato e anche il confrontato con alcuni algoritmi di clustering spettrale. Infine studiamo un criterio per decidere in quanti cluster partizionare l'insieme di partenza.

3.2.1 Clustering gerarchico

Una classe di algoritmi che prescinde dall'embedding in spazi vettoriali è quella *gerarchica*, in particolare *agglomerativa*[17]: “dal basso verso l'alto”, in cui ogni sequenza è inizialmente considerata come un cluster banale e coppie di cluster vengono aggregate man mano che si sale lungo la gerarchia.

In generale è sufficiente avere un insieme di elementi da clusterizzare e una matrice di dissimilarità che serve per indicare come unire i cluster. La matrice di dissimilarità deve essere simmetrica a entrate positive – criteri esattamente soddisfatti dalla matrice di distanze derivate dalla metrica di Rohlin. Al primo passo definiamo n cluster banali, contenenti un solo elemento; in questo modo possiamo parlare sempre di cluster, e non distinguere fasi intermedie con cluster e elementi ancora non aggiunti. Ad ogni passo i due cluster che hanno distanza *minima* tra loro vengono aggregati in un nuovo cluster binario. Rimane da definire la distanza tra due generici cluster \mathcal{A} e \mathcal{B} :

- Single linkage: la distanza tra cluster è definita come la minima distanza tra le coppie di elementi appartenenti ai due cluster.

$$d(\mathcal{A}, \mathcal{B}) = \min \{ d(x, y) : x \in \mathcal{A}, y \in \mathcal{B} \}$$

Non è una buona scelta in quanto la probabilità di aggregazione cresce con le dimensioni dei cluster presenti, i primi cluster formati quindi continuano a crescere per aggiunta di nuovi elementi – si forma un'albero molto sbilanciato in cui tutte le sequenze sono aggiunte gradualmente al primo cluster

- UPGMA, o average linkage: si prende la distanza media tra le coppie

$$d(\mathcal{A}, \mathcal{B}) = \frac{1}{|\mathcal{A}| \cdot |\mathcal{B}|} \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} d(x, y)$$

con un comportamento intermedio tra single e complete

- Complete linkage: la distanza è la massima tra tutte le coppie degli elementi

$$d(\mathcal{A}, \mathcal{B}) = \max \{ d(x, y) : x \in \mathcal{A}, y \in \mathcal{B} \}$$

i cluster tendono a formarsi di dimensioni simili e ben distinti tra di loro, per questo motivo è il criterio da noi scelto

Dopo l'ultimo passaggio tutti gli elementi si trovano in un unico cluster, che contiene un albero di sottocluster – rappresentabili con un dendrogramma, in cui l'altezza di ogni ramo è pari alla distanza tra le foglie. Partendo dall'alto e tagliando all'altezza minima possibile si possono ottenere p cluster dall'insieme iniziale (figura 3.2.1), con un algoritmo completamente generico e dipendente solo dalla matrice delle distanze.

Gli algoritmi gerarchici cercano una soluzione approssimata del problema del clustering in maniera *greedy* e hanno lo svantaggio che gli elementi una volta associati in un cluster non sono più spostati – un'associazione subottimale in partenza si propaga senza correzioni. Da qui la necessità del confronto con un altro metodo per avere la conferma che i cluster visti non sono stati introdotti dalla scelta del linkage completo. Il vantaggio della greediness è computazionale, con un andamento in tempo $\mathcal{O}(N^2)$ [18, 19] e l'utilizzo in memoria solo della matrice delle distanze. È possibile fare in un secondo decina di clusterizzazioni utilizzando i pacchetti inclusi in Matlab (vi sono reimplementazioni open altrettanto performanti in Python e R). In codice Matlab, dalla matrice *distanza*, con p clusters richiesti si ottiene il vettore *c* di etichette intere comprese tra 1 e p :

```
Z = linkage(squareform(distanza), 'complete');
c = cluster(Z, 'maxclust', p);
```

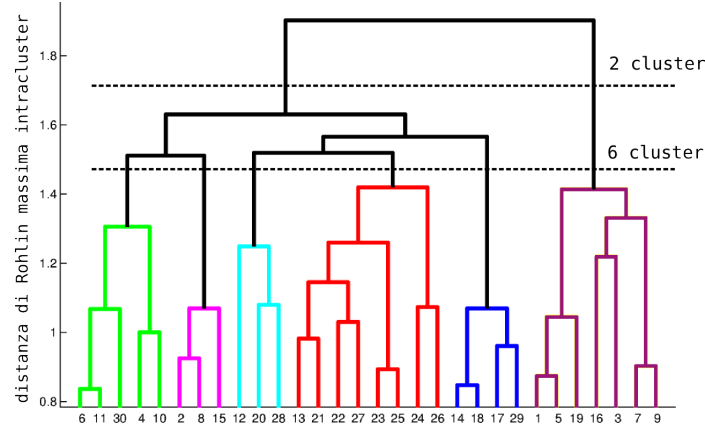



Figura 3.2.1: Dendrogramma costruito dai dati sulle distanze delle sequenze del H3N2 (mostrata solo parte superiore e non tutte le foglie). Sulle ordinate la distanza massima amessa dal complete linkage per accoppiare sottocluster; si nota come salendo, cluster diversi si trovano uniti fino a che non si raggiunge la massima distanza e l'accorpamento termina. Tagliando ad altezze opportune, abbiamo due esempi: un taglio che forma due cluster, un taglio che corrisponde ai 6 cluster colorati.

3.2.2 Clustering spettrale

Sotto questo nome cade un'ampia gamma di possibili metodi, basati tutti sugli autovettori di un laplaciano discreto opportunamente definito. Per un ottima review si veda [20]. Gli algoritmi si dividono grosso modo in due categorie:

- Riduzione dimensionale a partire dagli autovettori del laplaciano
- Taglio del grafo in sottopartizioni

Le sequenze che vogliamo clusterizzare formano i vertici di un grafo completo pesato G , in cui ogni sequenza è collegata ad altre con un peso determinato dalla similarità, che possiamo ottenere dalla matrice delle distanze tramite la funzione di similarità gaussiana

$$A(i, j) = \begin{cases} \exp\left(-\frac{d(i, j)}{2\sigma}\right) & \text{se } d(i, j) \neq 0 \\ 0 & \text{se } d(i, j) = 0 \end{cases}$$

definita a diagonale nulla, con parametro σ che la distanza tra vicini. Per ogni vertice i indichiamo il grado, fornito dalla somma dei pesi con tutti gli altri vertici

$$d_i = \sum_j A_{ij}$$

e la matrice diagonale ottenuta dai gradi così definiti con D .

Poiché il grafo è pesato con pesi simmetrici positivi $A(i, j) = A(j, i) \geq 0$, è completo e non diretto, possiamo definire il laplaciano nonnormalizzato come

$$L = D - A$$

Il laplaciano ha molte proprietà:

- L è simmetrico e definito positivo
- su ogni vettore $f \in \mathbb{R}^n$ agisce come una forma quadratica

$$\langle f, Lf \rangle = \sum_{i,j=1}^n A_{ij} (f_i - f_j)^2$$

- l'autovalore più piccolo è 0, con autovettore costante $(1, 1, \dots, 1)$
- la molteplicità dell'autovalore 0 è il numero di componenti connesse, ogni autovettore è l'equivalente discreto della funzione caratteristica sulle componenti connesse, con entrate 1 corrispondenti ai vertici che comprende e 0 altrove
- L ha autovalori reali nonnulli, $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$

Possiamo anche definire un laplaciano normalizzato come

$$L_{rw} := I - D^{-1}A$$

dove il pedice “rw” è dovuto alla similarità della definizione in problemi di random walk. Tutte le proprietà elencate per L valgono anche per L_{rw} invariate, sono solo riscritti gli autovalori. Indichiamo con “i primi p autovettori” gli autovettori corrispondenti agli autovalori *più piccoli* non nulli, ovvero scartando gli autovettori indicatori delle componenti connesse.

3.2.2.1 Riduzione dimensionale

L'algoritmo per il clustering di n elementi in p clusters, si basa sulla costruzione di L o L_{rw} ed è il seguente[21]:

1. In input la Matrice di similarità $A(i, j)$ e il numero p di cluster da costruire
2. Costruiamo la matrice dei gradi D , e il laplaciano $L_{rw} = I - D^{-1}A$
3. Calcolo dei primi p autovettori generalizzati u_1, \dots, u_p del problema agli autovalori $L_{rw} u = \lambda u$
4. Sia $U \in \mathbb{R}^{n \times p}$ la matrice contenente i vettori u_1, \dots, u_p come colonne.
5. Per ogni $i = 1, \dots, n$, sia $y_i \in \mathbb{R}^p$ il vettore corrispondente alla riga i -esima di U , ovvero ad ogni punto dell'insieme di partenza è stato fatto y_i
6. Troviamo ora i cluster C_1, \dots, C_p tramite l'algoritmo *kmeans* sui punti y di \mathbb{R}^p

Al punto 5 avviene la cosa più interessante, ad ogni elemento dell'insieme iniziale, caratterizzato solo dalle sue relazioni di distanza e similarità con gli altri elementi, è stato associato un punto nello spazio \mathbb{R}^p – la riduzione dimensionale permette di utilizzare, se il numero p è piccolo, *kmeans* con successo. Abbiamo trovato una rappresentazione funzionale per il clustering.

In figura 3.2.2 illustriamo il risultato su un solo autovettore – sono già evidenti alcuni cluster. Gli autovalori invece danno una indicazione di quanto sono connessi tra di loro i cluster: primi autovalori piccoli, prossimi allo zero, indicano un cluster quasi sconnesso; un grande *gap spettrale* indica invece una migliore coesione dei cluster (figura 3.2.3). La conta degli autovalori sotto un certo valore dà un'idea del numero di cluster ottimale.

3.2.2.2 Taglio del grafo

Si può procedere diversamente per partizionare un grafo in cluster: si può eseguire un taglio per dividere in due pezzi separati il grafo; si itera il procedimento di taglio su ogni sottografo così creato fino a raggiungere il numero desiderato di cluster. Bisogna cercare un criterio per ottenere un taglio ottimale – selezionare la componente con meno legami con il resto non è sufficiente, in quanto il più delle volte si sconnettono punti isolati, come si vede in figura 3.2.4.

Introduciamo un paio di definizioni: dati due sottoinsiemi di vertici A e B , introduciamo il peso dei legami relativi

$$W(A, B) = \sum_{i \in A, j \in B} A(i, j)$$

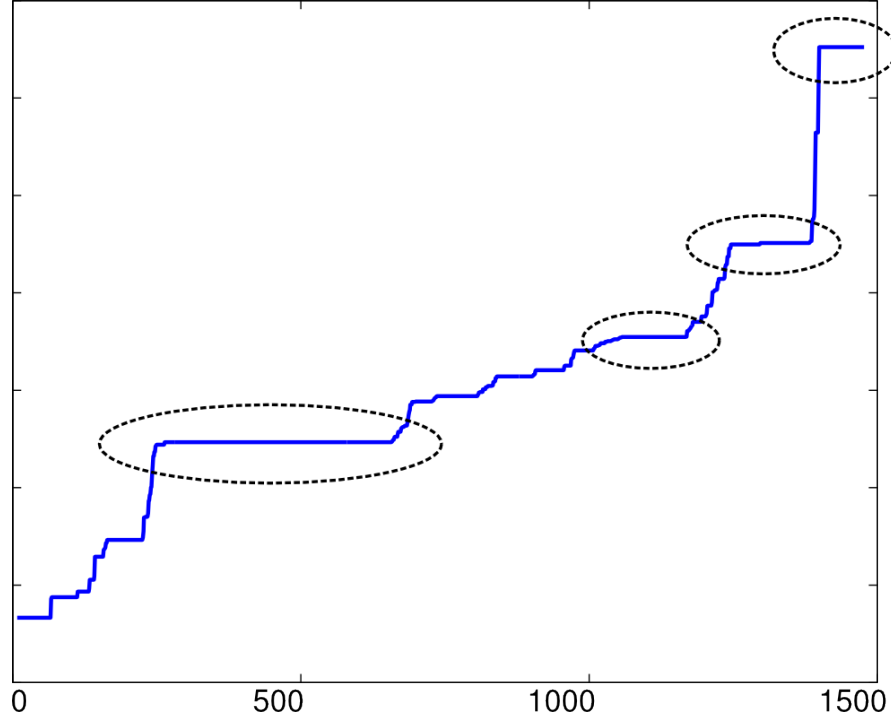


Figura 3.2.2: Andamento del secondo autovettore del laplaciano normalizzato. Il vettore è stato ordinato per valore, in modo che valori simili si trovassero vicini – ciò rende evidente la struttura a gradini, che grossomodo corrispondono ai cluster cercati, evidenziata in alcuni casi più evidenti.

Per ogni sottoinsieme del grafo A_k , con complemento \bar{A}_k , possiamo definire la cardinalità ed il volume

$$|A| = \text{numero dei vertici in } A$$

$$\text{vol}(A) = \sum_{i \in A} d_i$$

che pesano rispettivamente il numero di vertici e l'importanza per il resto del grafo, misurata in termini di peso dei link.

Indicando con A_1, \dots, A_p una possibile partizione del grafo il più semplice criterio di taglio è quello che trova direttamente la partizione A_1, \dots, A_p che minimizza il funzionale

$$\text{cut}(A_1, \dots, A_p) = \frac{1}{2} \sum_k W(A_k, \bar{A}_k)$$

ovvero riduce al minimo i legami tra le parti tagliate e il complementare. Tuttavia come si è visto in figura 3.2.4 porta ad un risultato non soddisfacente – bisogna i tagli normalizzare con il volume dei cluster formati, per evitare la formazione di cluster isolati. Ottimizziamo quindi il funzionale *normalized-cut*, o *ncut*:

$$\text{ncut}(A_1, \dots, A_p) = \frac{1}{2} \sum_k \frac{W(A_k, \bar{A}_k)}{\text{vol}(A_k)} = \sum_k \frac{\text{cut}(A_k, \bar{A}_k)}{\text{vol}(A_k)}$$

che rappresenta il giusto compromesso, normalizzando ogni cluster formato con il peso dell'insieme dei legami del cluster. Il criterio più popolare in letteratura è il cosiddetto *Rcut*, in cui al posto del volume, si normalizza con la cardinalità dei cluster formati

$$\text{Rcut}(A_1, \dots, A_p) = \frac{1}{2} \sum_k \frac{W(A_k, \bar{A}_k)}{|A_k|}$$

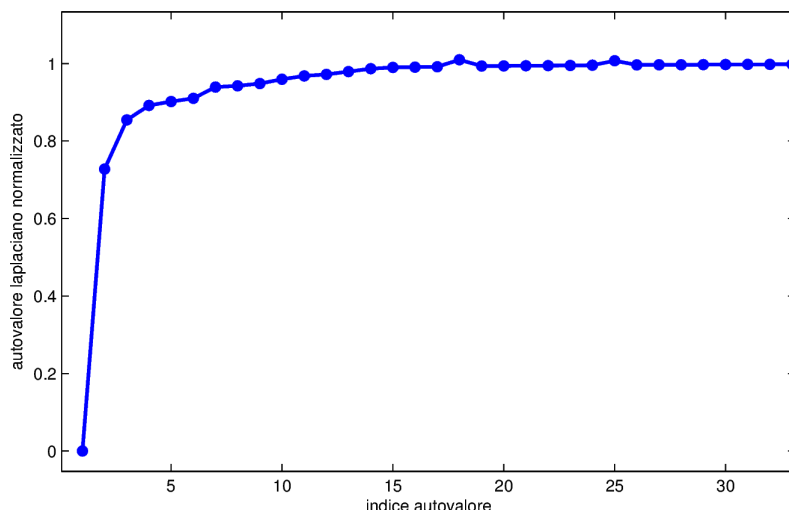


Figura 3.2.3: Primi autovalori del laplaciano normalizzato. Vi è un evidente gap tra l'autovalore zero e il primo nonnullo che significa una difficile separazione dei dati in cluster. Gli autovalori successivi sono molto vicini tra di loro, anche se si potrebbe riconoscere ad occhio che i primi 5 sono separati dai successivi, suggerendo un numero di cluster almeno pari a 5.

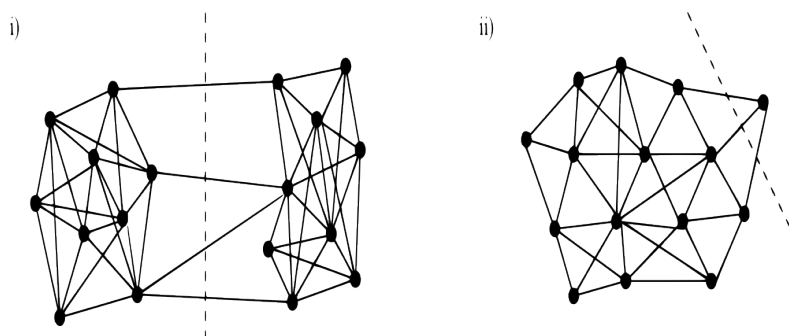


Figura 3.2.4: Criteri per il taglio del grafo. i) taglio "ottimale" ii) taglio che conta solo i legami rimossi. Illustrazione tratta da [22].

ma non è adatto ai nostri scopi, in quanto tende a selezionare cluster in cui il numero dei vertici (delle sequenze per noi) è quanto più simile – sappiamo invece che abbiamo molti dati recenti e pochi dell'inizio anni 90, per cui cluster "vecchi" avranno una cardinalità molte volte inferiore ai nuovi.

La soluzione esatta è nuovamente NP-completa e si cerca uno schema di approssimazione – si scopre[21][pag. 12, 20] che il modo migliore di approssimare le p partizioni è dato dall'utilizzo dell'algoritmo descritto in 3.2.2.1 a pagina 34. Tuttavia l'utilizzo di kmeans può portare in pratica a dei problemi per p grandi, come la mancata convergenza di kmeans, scelte povere dei centroidi iniziali, ecc.

In [22] si dimostra come l'applicazione ricorsiva, utilizzando il secondo autovalore del laplaciano per bipartire il grafo, è ottimale e ha *worst-case* garantito – la bipartizione del grafo è una clusterizzazione in una dimensione, caso semplice che non ha patologie. Chiaramente ad ogni iterazione bisogna costruire il laplaciano e ricalcolare il valore dell'autovettore su ogni sottografo considerato, un processo che può sembrare lungo, ma nel caso di grafi molto sparse è altamente ottimizzato[23] e rappresenta lo stato dell'arte per il clustering di preferenze commerciali (Amazon, ebay, ecc).

Per il calcolo utilizziamo il software creato per [24], che ottimizza un diverso funzionale, che tende nuovamente a equilibrare la cardinalità dei cluster, tramite una generalizza-

zione nonlineare del laplaciano

$$\langle f, L_k f \rangle = \sum_{i,j=1}^n A_{ij} |f_i - f_j|^k$$

Il software è comunque in grado di utilizzare il più semplice criterio di *ncut* quando $k=2$ e partizionare ricorsivamente il nostro insieme in maniera estremamente efficiente, generando anche tutti i clusters da intermedi 1 a p .

3.2.3 Numero ottimale di clusters

In tutti gli algoritmi di clustering esposti, è necessario stabilire *a priori* il numero dei cluster cercati. Questo presenta naturalmente uno svantaggio, visto che vogliamo scoprire la vera struttura dei cluster e non imporla. Cerchiamo quindi un criterio che codifichi l'aspettazione naïve nei confronti di un buon algoritmo, ovvero che i cluster “non siano né troppo grossi, né troppo spezzettati”.

Un modo per ottenere questo, è studiare a posteriori l'entropia dei cluster al variare del numero p di cluster imposti[15]. I cluster corrispondono a una partizione, per cui è possibile calcolarne l'entropia in funzione di p con la formula

$$H(p) = - \sum_{i=1}^n \frac{|A_i|}{N} \ln \frac{|A_i|}{N} \quad (3.2.1)$$

Per il clustering gerarchico, ci aspettiamo che per p molto piccolo vengano separati i cluster più diversi, corrispondenti ai primi rami del dendrogramma. Aumentando p , si arriva ad un punto corrispondente al “vero” numero di cluster, grossi e distinti. Crescendo p , ai cluster più grossi vengono rimossi gli “outliers”, piccoli in numero, ma distanti dal grosso del cluster. Aumentando sempre di più p , si arriva ad un punto in cui si spezza un grosso cluster, provocando un salto nell'entropia. L'entropia è per forza strettamente monotona rispetto a p , ma piccoli cluster con $\sim 1/N$ elementi hanno contribuito tendente a zero per N grande; una partizione con due cluster di $\frac{N}{2}$ siti ($H = \ln 2$), in cui uno viene spezzato a metà, vede aumentare l'entropia di

$$\Delta H = \frac{1}{2} \ln \frac{1}{2} - 2 \times \frac{1}{4} \ln \frac{1}{4} = \frac{\ln 2}{2} \simeq 50\%$$

Quindi “salti” nell'entropia indicano la rottura di un grosso cluster, mentre una crescita molto lenta con p indica l'allontanamento di elementi poco vicini agli altri. Il criterio da noi adottato è quindi la scelta di p come il punto in cui inizia una lenta crescita dell'entropia, lontano il più possibile da salti, ovvero la ricerca di un *plateau* dell'entropia.

L'applicabilità del criterio al clustering gerarchico è dovuta all'*effetto valanga* nel taglio del dendrogramma: vi sono p critici per i quali è necessario spezzare cluster grossi e p invece che causano solo minime modifiche. Il clustering spettrale non tende formare plateau al variare di p , poiché l'algoritmo tende a rompere i cluster in modo simile ad ogni passaggio $p \rightarrow p+1$.

Per il clustering gerarchico, l'applicazione è semplice e dà ottimi risultati, come si vede in figura 3.2.5 in alto a sinistra. La scelta della distanza di Rohlin ridotta per partizioni omogenee mette in evidenza un grosso plateau, che non è presente scegliendo altre distanze – testimonia ciò la stabilità dei cluster rispetto a piccole variazioni in p .

Non si trovano in letteratura altri criteri perfettamente applicabili al nostro caso: [25] studia il caso di 2-3 cluster, perfettamente visibili dal secondo autovettore del laplaciano normalizzato, mentre nel nostro caso (come si vede in figura 3.2.2 a pagina 35) sono riconoscibili solo alcuni raggruppamenti, ma non tutti. Un'alternativa possibile per il clustering divisivo tramite taglio del grafo, è fermarsi quando non è il criterio per quale cluster tagliare scende sotto una certa soglia – un'alternativa da esplorare in futuri lavori.

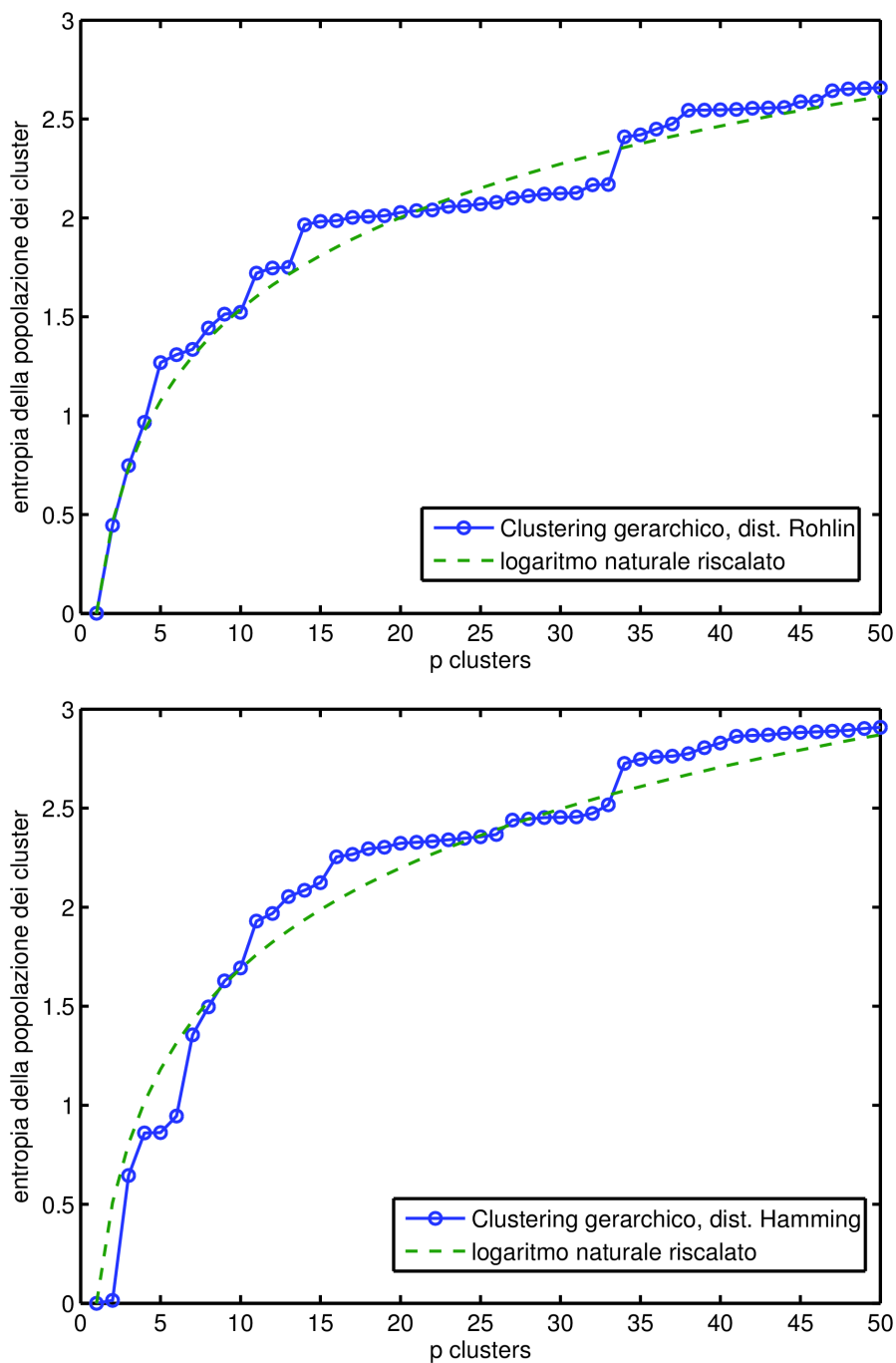


Figura 3.2.5: Andamento dell'entropia dei cluster generati con algoritmo gerarchico al variare di p . Le sequenze sono quelle del H3N2 senza deduplicazione, 1470 in totale. (A sinistra) la distanza è quella di Rohlin ridotta con partizioni omogenee. (A destra) la distanza di Hamming tra sequenze.

Usando la distanza di Rohlin si nota l'inizio del plateau per $p=14$, che prosegue fino a circa 30, indicando la stabilità dei cluster più grossi; con la distanza di Hamming l'andamento è meno regolare e il plateau molto meno marcato.

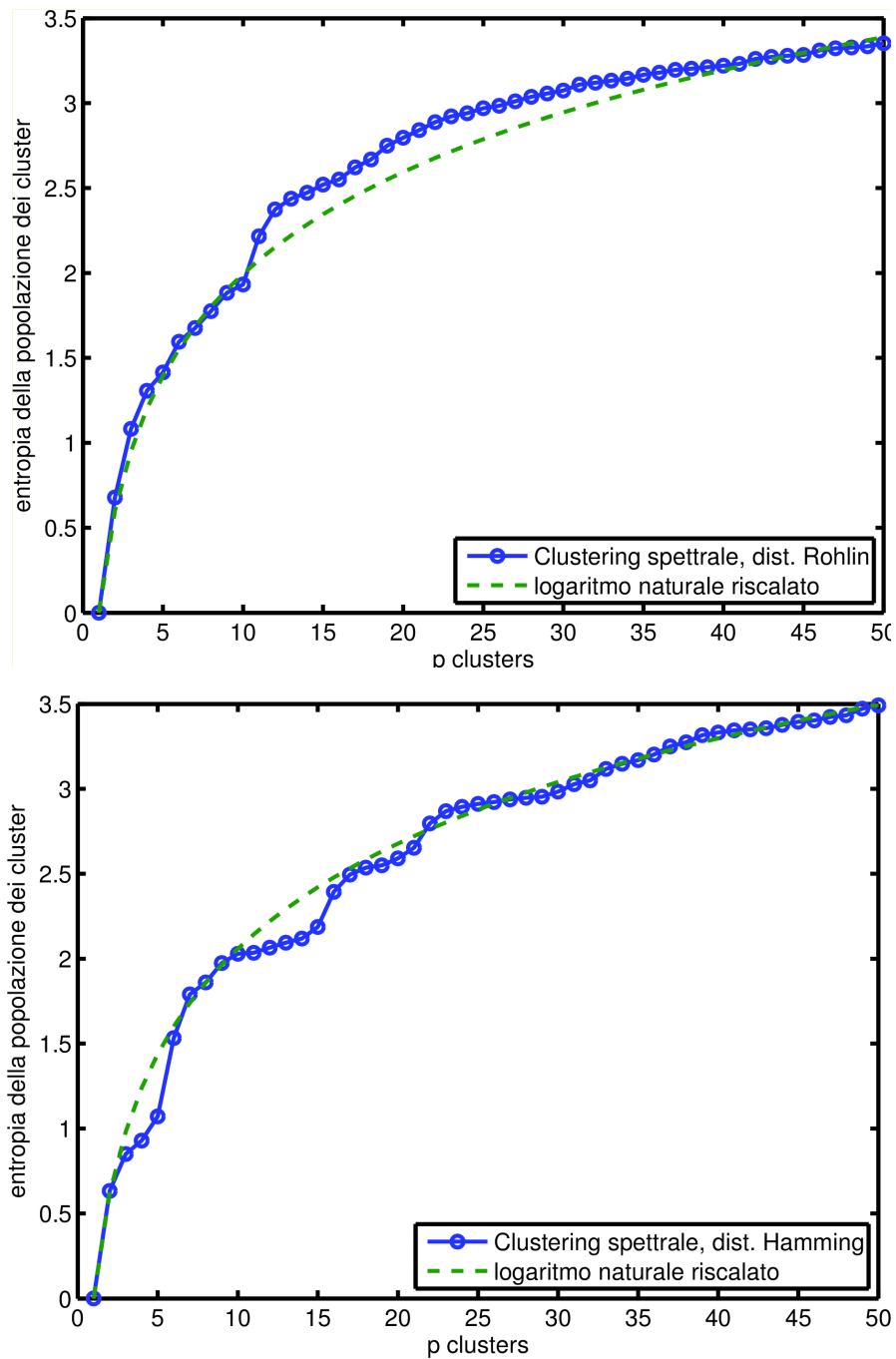


Figura 3.2.6: Andamento dell'entropia dei cluster generati con algoritmo spettrale al variare di p . Le sequenze sono quelle del H3N2 senza deduplicazione, 1470 in totale. (A sinistra) la distanza è quella di Rohlin ridotta con partizioni omogenee. (A destra) la distanza di Hamming tra sequenze. Nel caso del clustering spettrale il plateau è praticamente assente, per il diverso modo in cui vengono scelti i cluster ad ogni p .

Lunghezza	Conta
<560	2
564	3
565	2
566	1617
Totale	1624

a) HA

Lunghezza	Conta
<465	10
467-468	52
468	52
469	1510
Totale	1624

b) NA

Tabella 3.1: Distribuzione delle lunghezze delle sequenze di HA e NA come ottenute dal FluDB

3.3 Analisi preliminare delle sequenze

Vi sono numerosi database di materiale genetico dell'influenza, creati attorno al 2004 per organizzare in maniera sistematica i risultati del sequenziamento nei diversi paesi del mondo e rendere possibili rapide ricerche, confronti, analisi tra sequenze virali dell'influenza di ogni tipo. I database raccolgono materiale di svariate specie, dalle sequenze più antiche, comprendenti quelle del virus della pandemia del 1918 fino alle sequenze raccolte l'11 ottobre 2011 — non sono state ancora caricate sequenze della stagione influenzale corrente.

Vi sono svariati database, tutti permettenti la scelta del sottotipo di influenza, delle proteine da visualizzare, la dettagliata provenienza e data, tutti finanziati dal NIAID, con sequenze dal Influenza Genome Sequencing Project, GenBank e altri database non specializzati, rendendo semplice la ricerca, il confronto e l'analisi di molte sequenze diverse:

Influenza Research Database o IRD, noto anche come FluDB¹ [26]

Influenza Virus Resource organizzato dal NCBI ² [27]

Per rapportarsi ai risultati in [15] si è scelto le sequenze dal 1993 al presente, raccolte nel continente nord-americano, in cui la emoagglutinina è sequenziata completamente (al massimo in paio di “buchi”). Quando si è fatto il confronto con le sequenze della neuroaminidasi, sono state scelte le sequenze che hanno entrambe le proteine complete.

Rimane il problema di cosa fare per le sequenze con “buchi”: l'entropia e la distanza di Rohlin sono ben definite solo se le sequenze hanno la stessa lunghezza. Abbiamo quindi allineato le sequenze con un software apposito, Muscle[28], che inserisce sequenze di simboli '-' nei *gap* per riallineare le sequenze — prendiamo quindi un alfabeto esteso, che comprende i simboli aggiunti. Come si vede dalla tabella, il numero dei *gap* è piccolo se consideriamo la proteina HA, e anche nel caso della NA la maggior parte delle sequenze non ha *gap* e meno dell'1% ne ha più di 3.

L'alfabeto di base \mathbb{K} ha 20 simboli, ai quali si aggiungono 4 simboli che definiscono ambiguità, definiti in tabella 3.2. Il simbolo corrispondente ad amminoacido non identificato lo consideriamo distinto dal simbolo di allineamento '-'. Alla fine \mathbb{K} ha 25 simboli, il che non è un problema, in quanto il metodo di partizionamento distingue solo simboli uguali o diversi.

C'è da stabilire cosa fare con le sequenze ridondanti, ovvero le sequenze ottenute da campioni diversi, con la stessa sequenza di proteine HA e NA (se solo una è uguale, le consideriamo sempre distinte), che pure corrispondono a una notevole parte di quelle raccolte. Le caratteristiche qualitative del clustering non dipendono chiaramente dall'avere dei doppioni, poichè sono sicuramente messi nello stesso cluster; è utile avere le sequenze con doppioni anche per pesare i cluster e capire quali sono i più rilevanti — la numerosità indica la diffusione del particolare ceppo. La rimozione dei doppioni va comunque effettuata una volta scaricate le sequenze, i database online sono in grado di farla, ma

¹<http://www.fludb.org>

²<http://www.ncbi.nlm.nih.gov/genomes/FLU/FLU.html>

Amminoacidi ambigui	Simbolo
Asparagina o acido aspartico	B
Glutammina o acido glutamico	Z
Leucina o isoleucina	J
Non identificato	X oppure ''
Gap introdotto dall'alignment	–

Tabella 3.2: Tabella dei simboli aggiuntivi per l'alfabeto delle proteine

	Totale	Deduplicate	Ridondanti
H3N2 usate in [15]	1470	736	734 (50%)
H1N1 usate in [15]	2506	954	1552 (62%)
H3N2 dal IRD	1624	908	716 (44%)
H1N1 dal IRD	2908	1122	1784 (62%)

Tabella 3.3: Numero di sequenze considerate nei vari casi

identificano e rimuovono separatamente le sequenze corrispondenti a HA e NA, per cui non vi è una identificazione 1-1 tra i campioni delle due proteine e non si riesce a fare lo studio contemporaneamente.

Abbiamo quindi vari campioni di sequenze sia per l'influenza di tipo H3N2 che di tipo H1N1, a seconda delle scelte che si possono fare, elencati in tabella 3.3. Non vi sono particolari differenze nei risultati considerando le sequenze con o senza doppioni, ogni volta esplicheremo di quale dataset stiamo parlando.

Le sequenze sono molto simili, su una lunghezza media di 566 vi sono al massimo 50 lettere mutate tra le sequenze più distanti tra di loro – solo il 10% dei simboli quindi determina la distanza. Si vede dalla figura 3.3.1 come da uno sguardo rapido sia difficile individuare le mutazioni. Possiamo fare una mappa “caldo-freddo” per individuare i siti che mutano maggiormente, chiamati dai biologi *epitopi*.

La scelta del partizionamento omogeneo comporta un numero enorme di atomi, in generale i simboli tra un sito all'altro sono quasi sempre diversi. Si può studiare come si comporta il partizionamento con salto al variare del valore massimo di quest'ultimo (figura 3.3.3).

Da questi dati possiamo trarre le seguenti conclusioni, che influenzano i tipi di distanza e di clustering scelti:

- il numero di atomi è poco diverso da L , per l'atomo medio $\langle |A_k| \rangle \simeq 1$
- la maggior parte dei siti (e conseguentemente atomi) non cambia da una sequenza all'altra
- per qualunque tipo di partizionamento, la maggior parte degli atomi sarà identica, con una decina di differenze in media
- la riduzione avrà un effetto preponderante

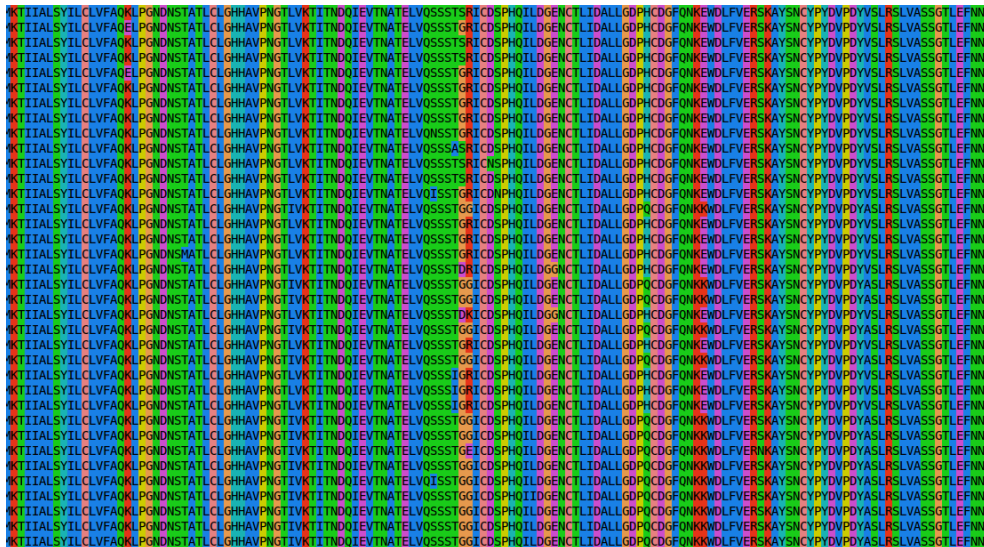


Figura 3.3.1: Vista della parte iniziale delle sequenze del H3N2, all'incirca 150 simboli

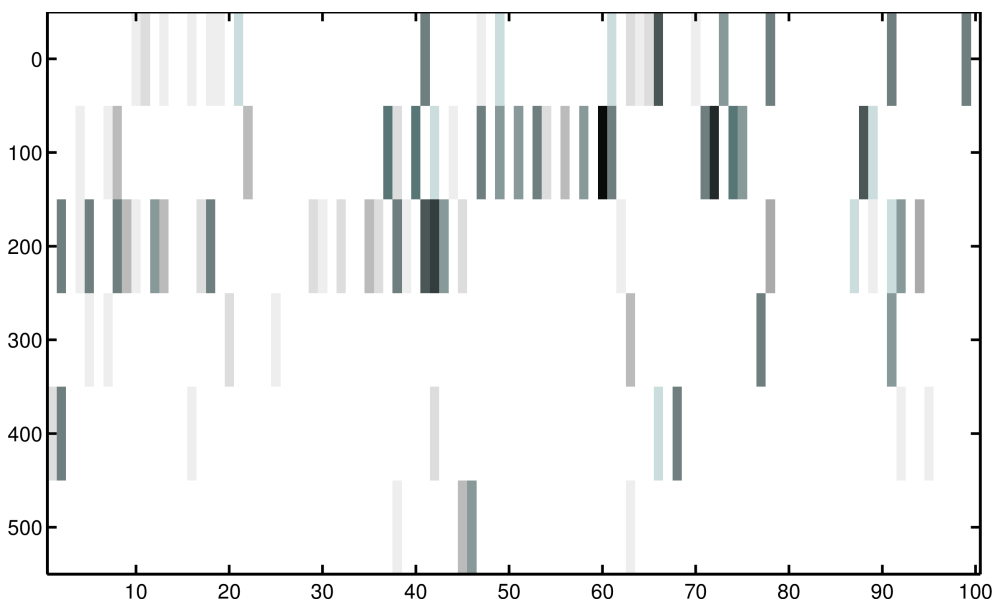


Figura 3.3.2: Mappa delle mutazioni per ogni sito. L'ordine è da sinistra a destra, dall'alto verso il basso. Poichè le sequenze sono lunghe 566, fino a 600 sono stati inseriti degli 0. In bianco i siti che non hanno avuto mutazioni, in scuro quelli con più cambiamenti – la quantità rappresentata è l'entropia dei simboli per sito, 0 dove sono tutti uguali e massima quando due o più amminoacidi hanno la stessa frequenza.

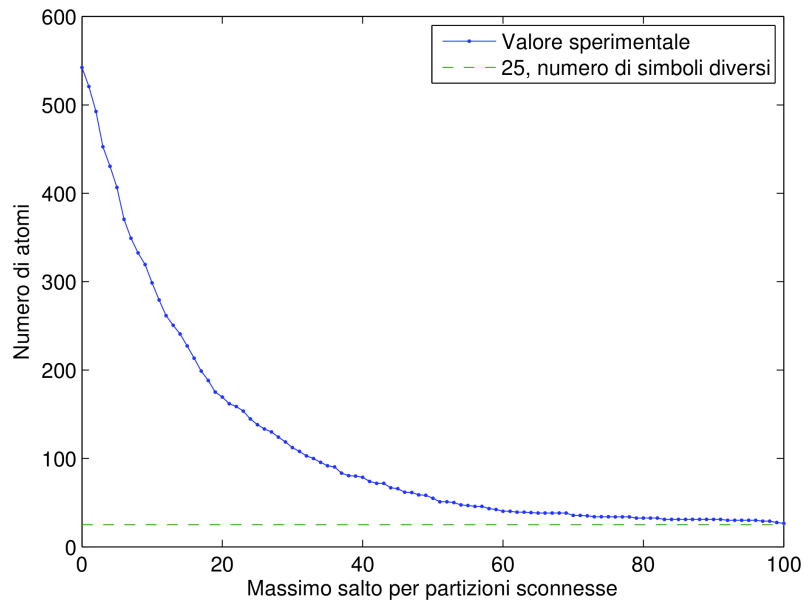


Figura 3.3.3: Andamento del numero di atomi medio per partizione di tipo generale non-connesso, al variare del salto massimo. Il numero parte dai 542 atomi medi per atomi connessi e si arriva a 25 atomi che raggruppano i tutti diversi simboli ammessi nelle sequenze.

3.4 Scelta della distanza

Abbiamo molte possibilità per definire una distanza usando i metodi introdotti nel primo capitolo, il numero delle diverse combinazioni cresce esponenzialmente e riassumiamo quelli calcolati ogni volta, come mostrato in tabella 3.4.

	Riduzione	Entropia	#
Distanza di Rohlin con partizionamento omogeneo	Senza	Shannon	1
		Topologica	2
	Con	Shannon	3
		Topologica	4
Distanza di Rohlin con tra partizioni "con salto"	Senza	Shannon	5
		Topologica	6
	Con	Shannon	7
		Topologica	8
Distanza Hamming	N/D	N/D	9

Tabella 3.4: Tabella distanze studiate

Esploriamo ora le caratteristiche di alcune distanze, mostrando quali scartare e perché, per illustrare nella prossima sezione la distanza più adatta al problema.

Osservazione matematica Possiamo dire quando la distanza di Rohlin è “piccola”: dalle proprietà generali della distanza sappiamo che

$$\begin{aligned} 0 &\leq d_R \\ 0 &\leq 2H(\alpha \vee \beta) - H(\alpha) - H(\beta) \quad \text{usando la definizione del prodotto} \\ 0 &\leq \frac{d_R}{H(\alpha) + H(\beta)} \leq 1 \quad \text{con } \alpha, \beta \text{ non entrambe nulle} \end{aligned}$$

per cui possiamo stimare che se la distanza è piccola o grossa confrontandola con l'entropia delle partizioni.

Ricordiamo che da due distanze è possibile definirne ∞^2 tramite combinazione lineare a coefficienti costanti positivi, ovvero se d_1 e d_2 sono distanze, lo è anche

$$d(\alpha, \beta) = a_1 d_1(\alpha, \beta) + a_2 d_2(\alpha, \beta) \quad a_1, a_2 > 0$$

per cui sarebbe possibile formare combinazioni lineari tra le distanze trovate, pesate opportunamente, per combinare i vantaggi e la evidenziazione delle caratteristiche cercate. L'enorme complessità aggiuntivo non ci ha indotto a perseguire questa strada.

3.4.1 Distanza di Hamming

La distanza di Hamming fornisce una prima ed utile indicazione con cui confrontare le misure successive. La semplice definizione ed il veloce calcolo, la rendono la distanza più utilizzata per la ricostruzione degli alberi filogenetici (eventualmente viene applicata una funzione della distanza). Tuttavia la sola considerazione delle mutazioni sui singoli siti non dà molta informazione sul cambiamento funzionale delle proteine coinvolte; la distanza risulta avere una componente proporzionale alla distanza temporale tra le sequenze – sappiamo che nel tempo l'antigenic drift accumula differenze, che rende difficile il riconoscimento dell'emergenza di un strain virale realmente nuovo con un numero ridotto di mutazioni antigeniche.

La distribuzione delle distanze di Hamming in figura 3.4.1 evidenzia che le mutazioni in quasi 20 anni di dati sulla struttura della molecola HA sono raramente più del 10%. Il clustering mette in evidenza l'andamento lineare con il tempo, senza la possibilità della distinzione di caratteristiche conservate nel tempo dai vari ceppi.

3.4.2 Distanza topologica

La distanza topologica tra due partizioni è calcolata sulla base del numero di atomi nel prodotto. Si utilizza la formula di Rohlin, con l'entropia definita dal logaritmo naturale del numero di atomi, che indichiamo con il simbolo n_α :

$$d_R^t(\alpha, \beta) = 2 \ln(n_{\alpha \vee \beta}) - \ln(n_\alpha) - \ln(n_\beta)$$

Focalizzando l'attenzione sul caso della distanza #2, ottenuta da partizionamento omogeneo senza riduzione, le partizioni ricavate dalle sequenze hanno tutte un numero molto simile di atomi e gli atomi stessi sono in stragrande maggioranza comuni ad entrambe le partizioni. Il prodotto quindi presenta la stessa parte comune dominante i pochi atomi non presenti in una delle due partizioni – il contributo dei frammenti, che viene misurato dalla distanza di Rohlin, è quasi trascurabile rispetto alla parte comune, e la distanza tende a zero. In più, come si vede in figura 3.4.2 a pagina 46, la distribuzione degli atomi nei prodotti $n_{\alpha \vee \beta}$ è estremamente simile alla distribuzione di n_α — il risultato è che la distanza topologica presenta una variabilità ridotta nello studio della distanza tra partizioni molto frammentate con una grandissima parte in comune. La distanza infatti è la differenza del logaritmo di due numeri molto grandi vicini e possiamo fare una stima della distanza media

$$\begin{aligned} \langle d_R^t \rangle &\sim \langle 2 \ln n_{\alpha \vee \beta} \rangle - 2 \langle \ln n_\alpha \rangle = 2 \ln(544.4) - 2 \ln(542.4) \\ &= 2 \ln \left(\frac{540}{540} \cdot \frac{1 + \frac{4.4}{540}}{1 + \frac{2.4}{540}} \right) \simeq 2 \cdot \frac{4.4 - 2.4}{540} = 0.007 \end{aligned}$$

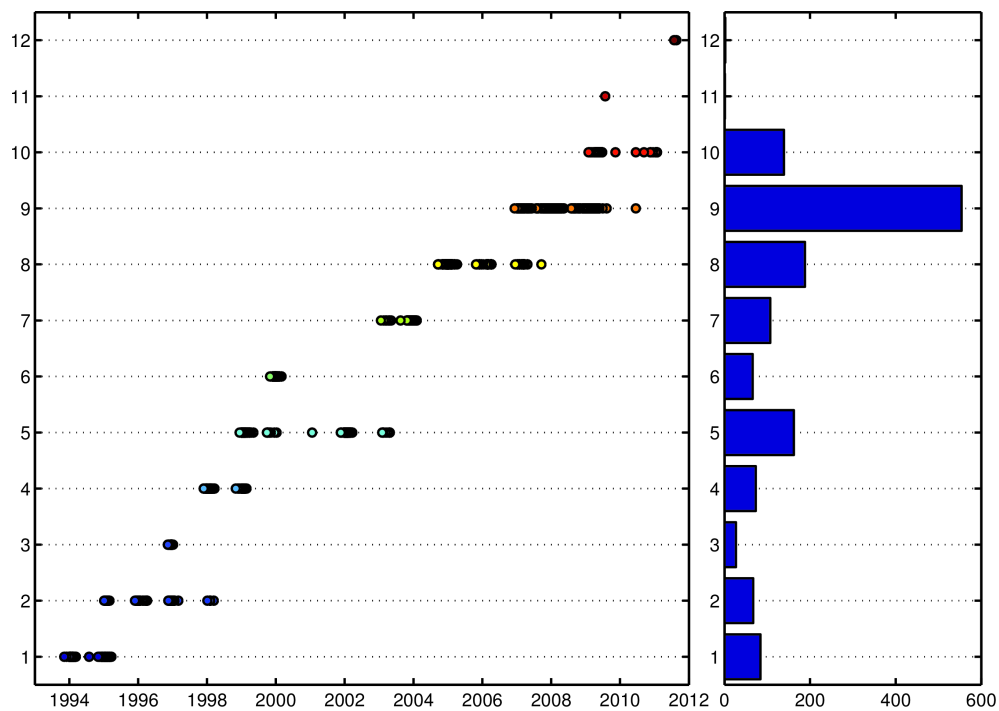
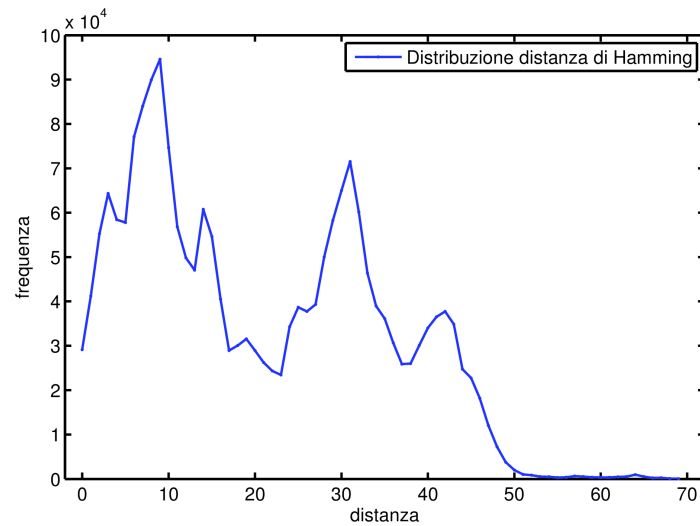


Figura 3.4.1: Dati calcolati tramite la distanza di Hamming sulle sequenze della HA del H3N2 con 1470 sequenze. (In alto) Distribuzione delle distanze. (In basso) Clustering, dal pronunciato carattere lineare del tempo.

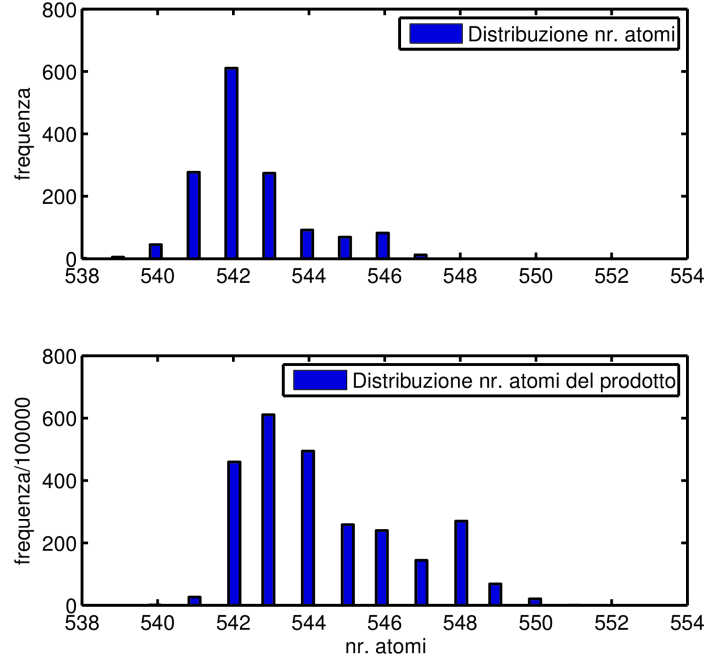


Figura 3.4.2: Distribuzione del numero degli atomi, in alto la distribuzione nelle sequenze, in basso nei prodotti. La scala degli istogrammi è diversa, in quanto a fronte di $N=1470$ partizioni ci sono $N^2/2$ prodotti.

contro i 0.008 misurati usando la vera distribuzione delle distanze, mostrata in figura 3.4.3. Per avere un'idea se la distanza è soppressa, calcoliamo

$$\frac{d_R^t}{2 \langle \ln n_a \rangle} \simeq \frac{7 \cdot 10^{-3}}{6.3} \simeq 10^{-4} \ll 1$$

che è molto vicino allo zero; il numero ridotto di combinazioni di $n_{\alpha\beta}$ implica che anche per le sequenze più distanti rispetto ad altre metriche, abbiamo una distanza topologica tendente a zero e senza grande variabilità.

La riduzione nelle partizioni connesse su sequenze lineari rimuove la parte comune, per cui la distribuzione del numero di atomi è praticamente invariata a meno di uno shift nell'intervallo $[0,10]$. Facendo una approssimata stima dell'incremento medio della distanza topologica grazie alla riduzione si ha

$$\langle R \rangle_{\text{stima}} = \frac{\langle d_R^t(\hat{\alpha}, \hat{\beta}) \rangle}{\langle d_R^t(\alpha, \beta) \rangle} = \frac{2 \ln \left(\frac{4.4}{2.4} \right)}{2 \cdot \frac{4.4-2.4}{540}} \simeq \ln(1.8) \frac{540}{2} \simeq 163$$

il dato sperimentale è

$$\langle R \rangle_{\text{misura}} = 183$$

per cui l'ipotesi del grossa parte comune che schiaccia la distanza topologica è corretta. Possiamo quindi essere tranquilli nel scartare questa distanza come poco significativa.

3.4.3 Distanza tra partizioni sconnesse

Abbiamo mostrato come il numero degli atomi cali molto, al crescere del numero di salti possibili. Vi sono però notevoli svantaggi, a partire da quello computazionale: le partizioni di tipo generico, non connesse, richiedono il doppio del tempo per il calcolo dei prodotti e ancora 4 volte tanto per la riduzione.

Presentiamo ora molte illustrazioni dei comportamenti delle distanze quando le partizioni sono lineari ma sconnesse. Non consideriamo la distanza topologica, in quanto

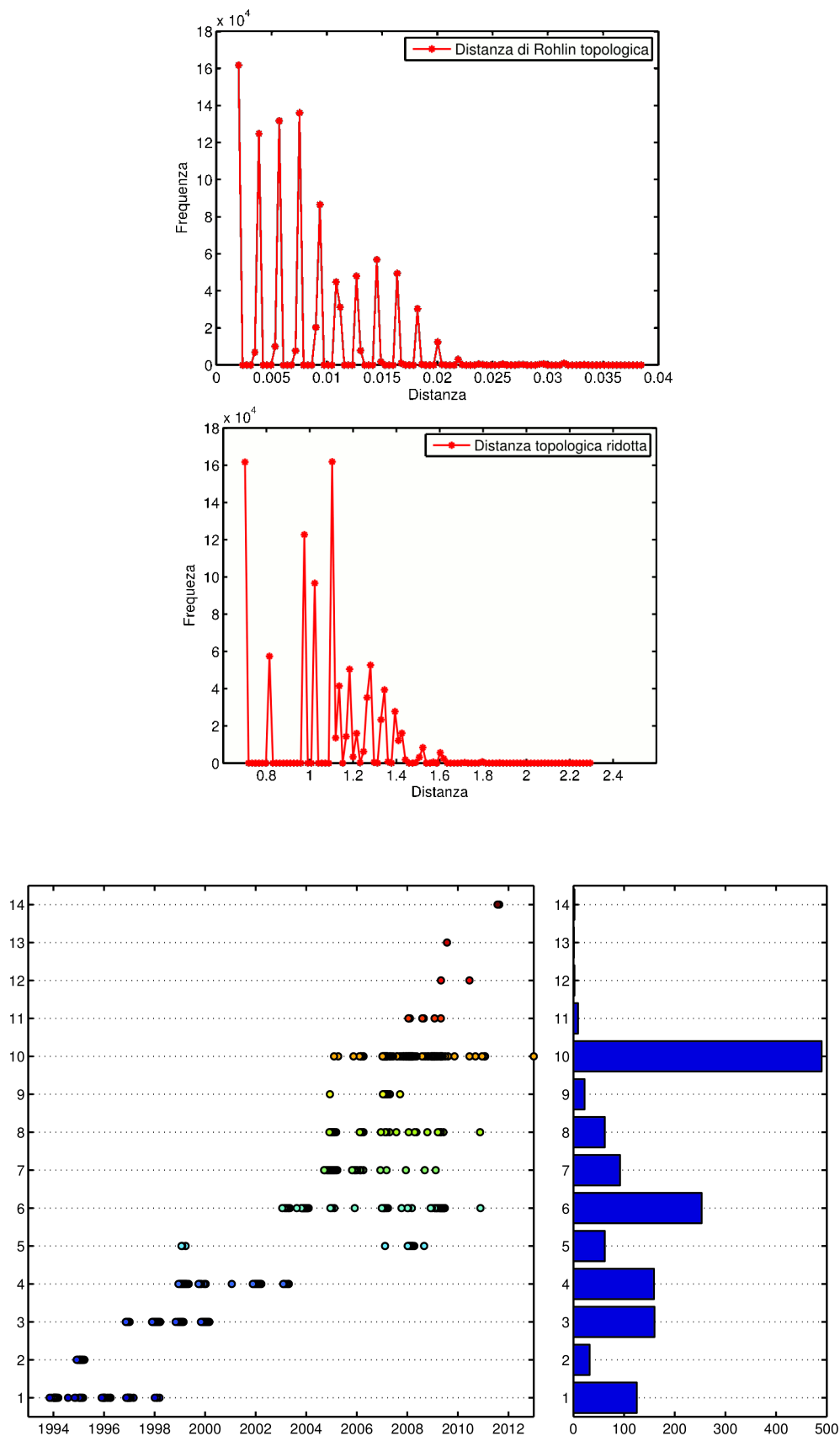


Figura 3.4.3: Dati sulla distanza topologica usando le sequenze del H3N2. (In alto) Distribuzione istogrammi delle distanze tra le sequenze, nel caso (a sinistra) senza riduzione e (a destra) con riduzione. (In basso) Clusterizzazione usando la distanza topologica nonridotta, $p=14$ per confrontarlo con i clustering successivi, in quanto l'andamento dell'entropia dei cluster non presentava plateau né altri punti notevoli.

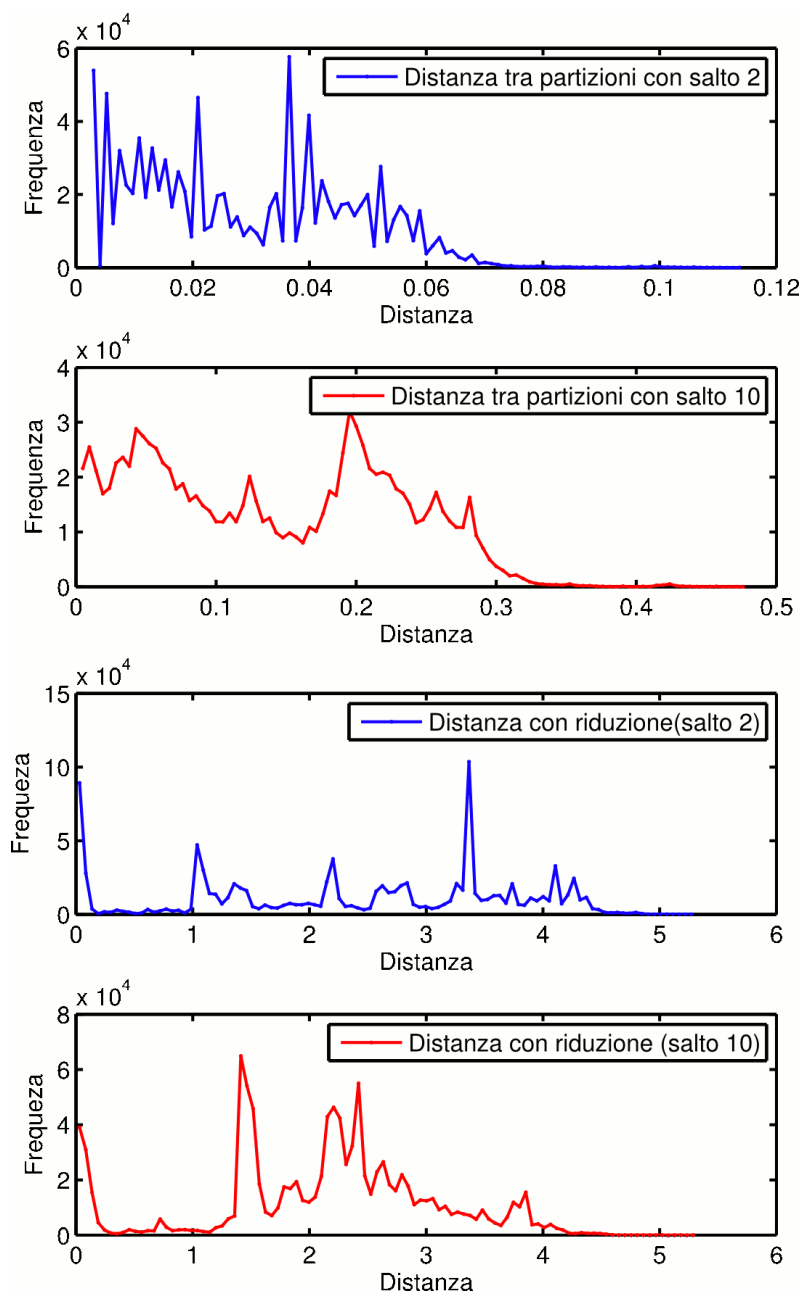


Figura 3.4.4: Confronto dei comportamenti delle distanze su partizioni “con salto”. Anche senza riduzione, la possibilità di avere atomi più estesi aumenta notevolmente le distanze tra le partizioni. La riduzione successivamente livella questa differenza, $R=88$ nel primo caso e $R=16$. Lo studio dettagliato del significato di queste distanze ridotte non è stato fatto, poiché si è visto che la clusterizzazione indotta è simile a quella ottenuta con Hamming – non presenta informazioni utili.

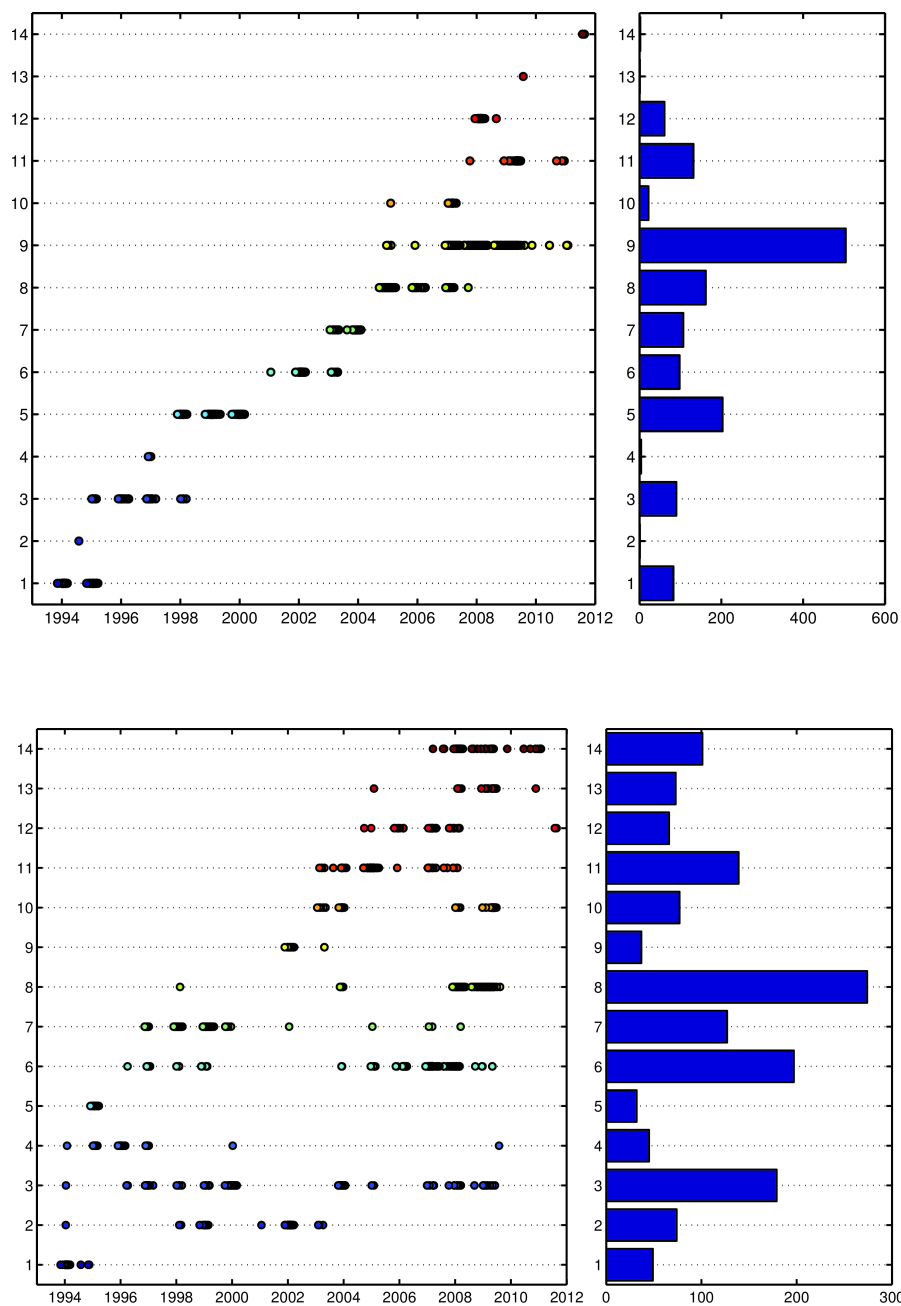


Figura 3.4.5: Illustrazione del clustering utilizzando due diverse distanze. (In alto) Distanza su partizioni che ammettono atomi con salto 2, con un andamento spiccatamente lineare nel tempo, simile a quanto accade con Hamming. (In basso) Salto 10, ma utilizzando le sequenze ridotte – non vi è più una linearità, ma i cluster sono diventati fin troppo estesi temporalmente, un altro comportamento che sembra non catturare la natura biologica del sistema. Questo tipo di partizionamento raggruppano maggiormente l'informazione presente nelle sequenze, fino a riottenere un comportamento simile alla distanza di Hamming, nel limite in cui si prenda n_s massimo.

La distanza tra partizioni ridotta con salto 2 e quella nonridotta con salto 10 assomigliano alla figura in alto, e non sono state incluse.

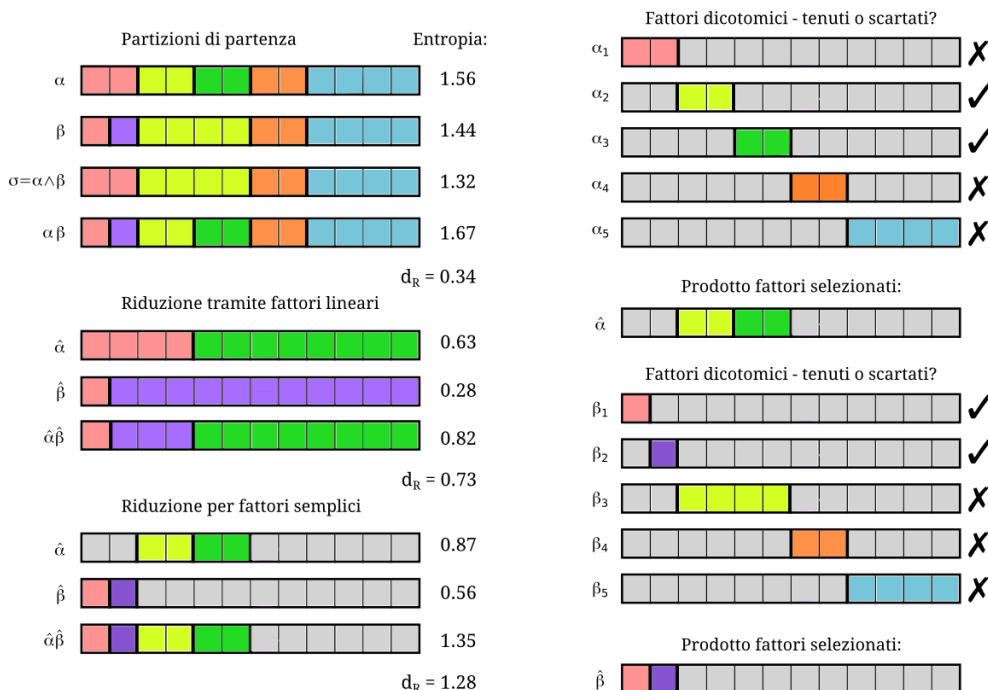


Figura 3.4.6: Dettaglio analisi dei due tipi di riduzione, con fattori lineari connessi e fattori dicotomici. Quando gli atomi non sono connessi vi è un solo tipo di riduzione possibile, il confronto in questo caso serve per chiarire l'effetto molto diverso del fattore comune σ . A destra di ogni partizione è indicata l'entropia di Shannon; alla fine di ogni blocco usando le entropie illustrate è possibile calcolare la distanza di Rohlin corrispondente.

ancor meno significativa. Tutti i casi trattati sono quelli della proteina HA del sottotipo H3N2, per poter fare confronti con gli altri.

La riduzione tramite fattori dicotomici universali non raggruppa gli atomi consecutivi comuni fino alla prima differenza tra le sequenze, ma raggruppa tutti gli atomi comuni in un unico grande atomo sconnesso, frammentato dai piccoli (e di poca rilevanza nell'entropia della partizione ridotta) superstiti. Questo fa sì che il prodotto avrà una entropia molto maggiore rispetto al prodotto di partizioni ridotte nel caso connesso. Il risultato è un grande fattore di amplificazione, ma una distanza risultante poco utile – troppo peso è dato alla parte comune nascondendo mutazioni che differenziano le sequenze. In figura 3.4.6 la spiegazione e dettaglio confronto, nel caso di due sequenze costruite ad hoc.

3.4.4 Proteine diverse, NA

Avevamo notato l'esistenza di altri dati, quelli riguardanti le sequenze della proteina antigenica *neuroaminidase*, perché tutte le misure sono state fatte invece sulla emoagglutina?

Un modo rapido per capire i motivi è verificare a che clustering porta una distanza eseguita sul NA – si vede come è una proteina molto più conservata rispetto ad HA, con minori mutazioni nel tempo, i cluster appaiono molto stabili e omogenei, indice di una scarsa differenza. Persino la misura fatta con la distanza di Hamming rivela che non vi sono rilevabili differenze. La proteina NA non è quindi adatta al tipo di misure che facciamo: il riconoscimento dei ceppi dovuti a antigenic drift tra stagioni influenzali vicine.

Le altre proteine del virus non sono antigeni, non sono cioè sottoposte a pressione evolutiva da parte del sistema immunitario e sono estremamente conservate rispetto alle proteine sul capsido, rendendole meno utili per l'identificazione, ma migliori target per futuri vaccini, gli studi sono in corso!

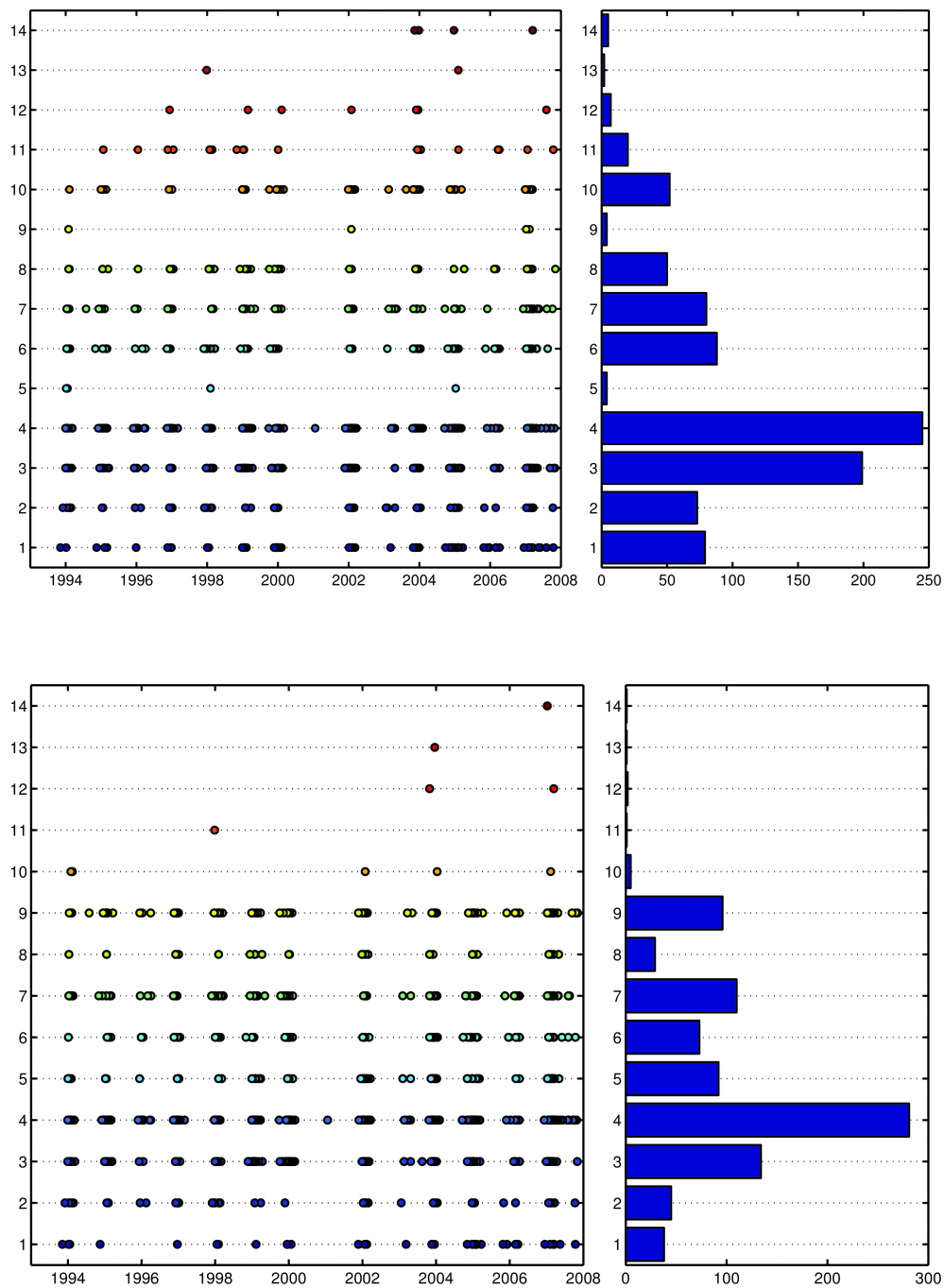


Figura 3.4.7: Clustering gerarchico agglomerativo basato sulle sequenze deduplicate della proteina NA per virus del sottotipo H3N2. Un'analoga misura per il sottotipo H1N1 non mostra alcuna differenza, per cui non replichiamo i risultati. (In alto) Distanza di Rohlin ridotta con fattori dicotomici lineari (In basso) Distanza di Hamming.

I risultati sono molto ambigui: innanzitutto abbiamo forzato $p=14$, anche se la clusterizzazione naturale indica $p \leq 5$, per poter fare il confronto con i dati della prossima sezione e per cercare di mostrare una qualunque eventuale struttura interna ai cluster. Non compare nulla di ciò, addirittura sembra che la distanza sia troppo poca per riempire i 14 cluster richiesti nel caso di Hamming. Le sequenze non tendono a clusterizzarsi né temporalmente, né in ceppi – la proteina NA è molto più conservata, andrebbe affrontata con un modo e una misura differenti.

3.5 Conclusione

Con un'attenta analisi di tutte le variazioni, abbiamo scartato 7 distanze delle 9 in tabella 3.4 a pagina 43. Siamo arrivati alla conclusione che la distanza migliore, che mette in risalto le caratteristiche delle sequenze, ha le seguenti caratteristiche:

- distanze calcolate tutte sulla proteina HA solamente
- distanza di Rohlin definita con l'entropia di Shannon
- partizioni lineari connesse
- riduzione con fattori lineari

Vediamo ora di motivare l'ultimo punto: la scelta dei fattori lineari connessi è in linea con quanto illustrato nelle figure precedenti. L'utilizzo della riduzione con fattori dicotomici ha avuto risultati indesiderati, perché i fattori lineari dovrebbero essere migliori?

Senza riduzione, visto che quasi tutti gli atomi sono banali, una mutazione ha in media scarso effetto – da un atomo con un solo elemento, si ha nuovamente un atomo con un solo elemento, le partizioni sono immutate e la distanza è zero. La distanza di Rohlin non ridotta ha quindi limitato potere risolutivo tra le effettive differenze dei ceppi virali.

La riduzione con fattori lineari rimuove i confini degli atomi in comune, che sono la maggioranza. Lascia quindi i bordi corrispondenti alle mutazioni, formando atomi pari al segmento tra una mutazione e la successiva: poichè l'entropia è metrica, il peso dipende dalla lunghezza di questi segmenti. Nei tratti in cui le mutazioni sono frequenti il peso è naturalmente più piccolo, mentre dove sono rare una mutazione è premiata. Questo particolare sistema di pesi, che entra nelle distanze tramite i prodotti tra partizioni, evidenzia le differenze in tratti in cui vi è minore entropia e sembra quindi mettere in mostra al meglio le caratteristiche biologiche delle sequenze, senza mai però richiedere una conoscenza dei meccanismi in gioco.

Mostriamo ora i risultati ottenuti usando questa distanza, in linea con quanto esposto in [15]: abbiamo confermato la scelta della distanza di Rohlin ridotta più semplice. Il numero di cluster ottimale $p=14$ è stato determinato in figura 3.2.5 a pagina 38 in alto, lo useremo per tutti i clustering qui mostrati. I dati provengono tutti dalle 1470 sequenze del H3N2 da fonte precedentemente citata.

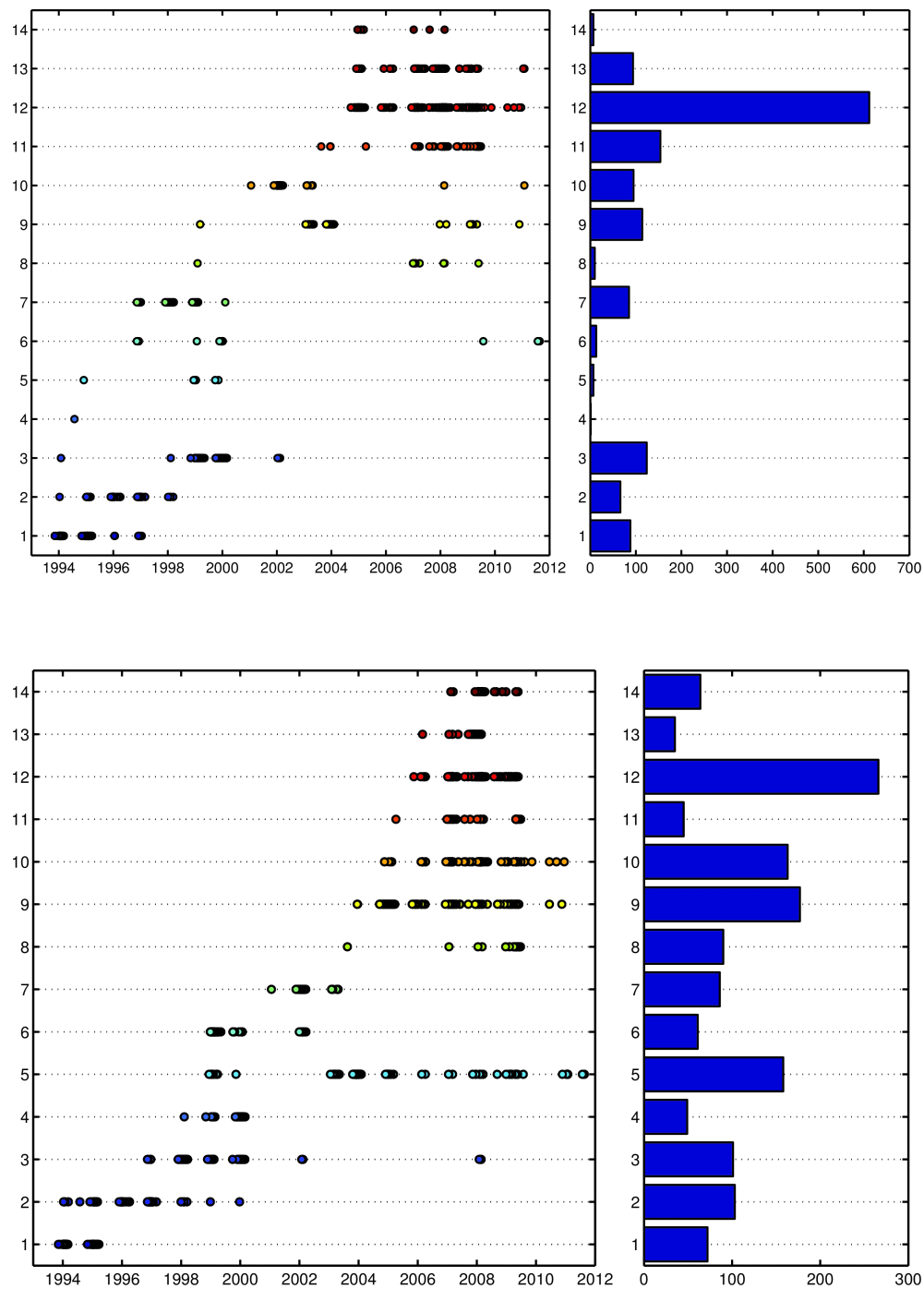


Figura 3.5.1: Clusters ottenuti con la distanza di Rohlin ridotta, usando fattori dicotomici lineari connessi. (In alto) Clustering gerarchico agglomerativo, tagliato a $p=14$. (In basso) Clustering spettrale divisivo, fermato al 13 passaggio. Si nota la forte somiglianza qualitativa dei due clusterings, ottenuti in modo completamente diverso. Questo rafforza la convinzione dell'effettiva importanza della distanza utilizzata per la classificazione delle sequenze di proteine HA.

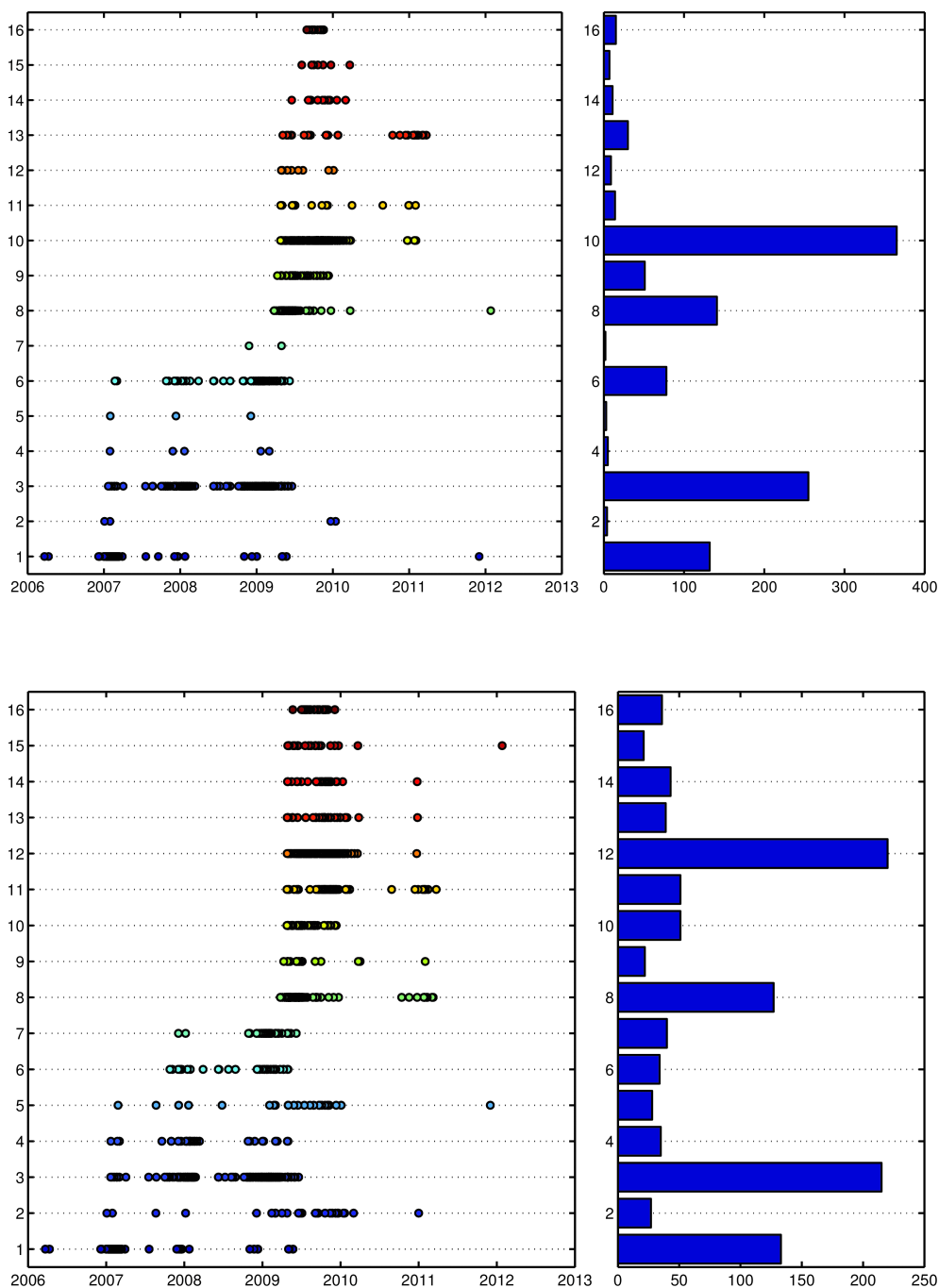


Figura 3.5.2: Utilizziamo stavolta il dataset del virus del sottotipo H1N1, con sequenze deduplicate. (In alto) Clustering gerarchico agglomerativo, tagliato a $p=16$. (In basso) Clustering spettrale divisivo, fermato al 15 passaggio. La clusterizzazione è decisamente compatibile e mostra il sorgere di almeno 8 ceppi distinti di H1N1 nel 2009, anno in cui un antigenic shift portò al diffondersi nell'uomo di ceppi novelli, prima diffusi solo in uccelli, in competizione tra loro. Si vede come all'avanzare del tempo sono rimasti solo due ceppi dominanti – lo shift ha portato una notevole diversità genetica che ha dovuto competere per la sopravvivenza.

Bibliografia

- [1] P. Billingsley, Ergodic theory and information, Wiley, New York, 1965. 1.1, 1.1
- [2] E. Agliari, M. Casartelli, E. Vivo, Metric characterization of cluster dynamics on the sierpinski gasket, Journal of Statistical Mechanics: Theory and Experiment 2010 (2010) P09002. 1.1
- [3] J. Munkres, Topology, Prentice Hall, Inc, Upper Saddle River, NJ, 2000. 1.1.1
- [4] J. Hoshen, R. Kopelman, Percolation and cluster distribution. i. cluster multiple labeling technique and critical concentration algorithm, Physical Review B 14 (8) (1976) 3438–3445. 2.2.2
- [5] M. Newman, R. Ziff, Fast monte carlo algorithm for site or bond percolation, Physical Review E 64 (1) (2001) 016706. 2.2.2
- [6] M. Mitzenmacher, S. Vadhan, Why simple hash functions work: exploiting the entropy in a data stream, in: Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms, SODA '08, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2008, pp. 746–755. 2.2.4
- [7] D. E. Knuth, Art of Computer Programming, Volume 3: Sorting and Searching, 2nd Edition, Vol. 3, Addison-Wesley Professional, 1998. 2.2.4
- [8] A. Partow, General purpose hash function algorithms.
URL <http://www.partow.net/programming/hashfunctions/index.html> 2.2.4
- [9] A. Appleby, Murmurhash 3.
URL <http://code.google.com/p/smhasher/> 2.2.4
- [10] N. M. Bouvier, P. Palese, The biology of influenza viruses., Vaccine 26 Suppl 4 (2008) D49–D53. 3.1
- [11] D. Goodsell, Hemagglutinin, pdb molecule of the month.
URL <http://www.rcsb.org/pdb/101/motm.do?momID=76> 3.1.1
- [12] F. Carrat, A. Flahault, Influenza vaccine: the challenge of antigenic drift., Vaccine 25 (39-40) (2007) 6852–6862. 3.1.2
- [13] M. F. Boni, Vaccination and antigenic drift in influenza., Vaccine 26 Suppl 3 (2008) C8–14. 3.1.2
- [14] WHO, Influenza vaccine viruses and reagents.
URL <http://www.who.int/influenza/vaccines/virus/en/> 3.1.2
- [15] R. Burioni, R. Scalco, M. Casartelli, Rohlin distance and the evolution of influenza a virus: Weak attractors and precursors, PloS one 6 (12) (2011) e27924. 3.1.2, 3.2.3, 3.3, ??, ??, 3.5

- [16] T. Hastie, R. Tibshirani, J. Friedman, The elements of statistical learning - data mining, inference, and prediction, 2nd Edition, Springer, Berlin, Heidelberg, 2008.
URL <http://www-stat.stanford.edu/~tibs/ElemStatLearn/> 3.2
- [17] G. Gan, C. Ma, J. Wu, Data clustering - theory, algorithms, and applications, illustrated edition Edition, SIAM, Society for Industrial and Applied Mathematics, 2007. 3.2.1
- [18] W. Press, S. Teukolsky, W. Vetterling, B. Flannery, Numerical Recipes 3rd Edition: The Art of Scientific Computing, Cambridge University Press, 2007. 3.2.1
- [19] R. Sibson, Slink: an optimally efficient algorithm for the single-link cluster method, The Computer Journal 16 (1) (1973) 30–34. 3.2.1
- [20] U. Von Luxburg, A tutorial on spectral clustering, Statistics and Computing 17 (4) (2007) 395–416. 3.2.2, 3.2.2.2
- [21] J. Shi, J. Malik, Normalized cuts and image segmentation, Pattern Analysis and Machine Intelligence, IEEE Transactions on 22 (8) (2000) 888–905. 3.2.2.1, 3.2.2.2
- [22] R. Kannan, S. Vempala, A. Vetta, On clusterings: Good, bad and spectral, Journal of the ACM (JACM) 51 (3) (2004) 497–515. 3.2.4, 3.2.2.2
- [23] D. Cheng, R. Kannan, S. Vempala, G. Wang, A divide-and-merge methodology for clustering, ACM Transactions on Database Systems (TODS) 31 (4) (2006) 1499–1525. 3.2.2.2
- [24] M. Hein, T. Bühler, An inverse power method for nonlinear eigenproblems with applications in 1-spectral clustering and sparse pca, Arxiv preprint arXiv:1012.0774.
URL <http://www.ml.uni-saarland.de/code/pSpectralClustering/pSpectralClustering.html> 3.2.2.2
- [25] C. Kim, M. Cheon, M. Kang, I. Chang, A simple and exact laplacian clustering of complex networking phenomena: Application to gene expression profiles, Proceedings of the National Academy of Sciences 105 (11) (2008) 4083.
URL <http://www.pnas.org/content/105/11/4083.full> 3.2.3
- [26] R. B. Squires, J. Noronha, V. Hunt, A. Garcia-Sastre, C. Macken, N. Baumgarth, D. Suarez, B. E. Pickett, Y. Zhang, C. N. Larsen, A. Ramsey, L. Zhou, S. Zaremba, S. Kumar, J. Deitrich, E. Klem, R. H. Scheuermann, Influenza research database: an integrated bioinformatics resource for influenza research and surveillance., Influenza Other Respi Viruses. 3.3
- [27] Y. Bao, P. Bolotov, D. Dernovoy, B. Kiryutin, L. Zaslavsky, T. Tatusova, J. Ostell, D. Lipman, The influenza virus resource at the national center for biotechnology information., J Virol 82 (2) (2008) 596–601, nCBI flu database. 3.3
- [28] R. C. Edgar, Muscle: multiple sequence alignment with high accuracy and high throughput., Nucleic Acids Res 32 (5) (2004) 1792–1797. 3.3