

Multiple Sequence Alignment (MSA) di sequenze SARS-Cov2

Silva Edoardo 816560, Marchetti Davide 815990

21/05/2020

Abstract

Usando i sequenziamenti genomici del virus denominato Covid-19, reperibili sul sito della banca genetica <https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/> e <https://www.gisaid.org/>, abbiamo allineato un insieme di sequenze relativi a paesi mediorientali:

- Iran: MT320891.2, MT281530.2, hCov-19/Iran/KHGRC-1.1-IPI-8206/2020|EPI_ISL_447469|2020-03-09, hCoV-19/Iran/HGRC-2-2162/2020|EPI_ISL_437512|2020-03-26.
- Israele: MT276598.1, MT276597.1, hCoV-19/Israel/130710062/2020|EPI_ISL_447469|2020-04-14.
- Pakistan: MT262993.1, MT240479.1, hCoV-19/Pakistan/Gilgit1/2020|EPI_ISL_417444|2020-03-04.
- Turchia: MT327745.1, hCoV-19/Turkey/HSGM-10232/2020|EPI_ISL_437334|2020-03-24, hCoV-19/Turkey/HSGM-1027/2020|EPI_ISL_437317|2020-03-27.

E allineati grazie ai tools Clustal Omega e MUSCLE , reperibili al sito <https://www.ebi.ac.uk/Tools/msa/>; al fine di ottenere le differenze con la sequenza di riferimento ottenuta su un campione di Whuan NC_045512.2.

Il lavoro consiste in identificare le differenze tra le sequenze e tra gli allineamenti dei 2 tool sulle medesime sequenze.

1 Descrizione

L'analisi filogenetica serve a ricostruire la storia delle mutazioni, ossia come specie antiche si sono evolute in quelle moderne.

I geni sono composti da sequenze di Acido Desossiribonucleico (DNA) (o RNA in caso di alcuni virus), composto da uno zucchero (desossiribosio) che unisce un gruppo fosfato ad una base azotata per comporre nucleotidi (i nucleotidi sono unità ripetitive costitutive degli acidi nucleici) che compongono il genoma.

Un GENE non è nient' altro che una particolare sequenza di DNA, che codifica l' informazione in un linguaggio a quattro lettere, nel quale ogni lettera è rappresentata da una base.

Il genotipo di un individuo è dato dal suo corredo genetico. Il fenotipo, invece, è l'insieme dei caratteri che l'individuo manifesta: dipende dal suo genotipo, dalle interazioni fra geni e anche da fattori esterni.

Il DNA viene utilizzato:

1. trascrivendolo in RNA
2. rendere l' RNA stabile aggiungendo ai limiti 7-metilguanosina (7mGTP)
3. trascritto il risultato precedente in amminoacidi che creeranno proteine.

2 Codice

Il codice è in python ed esegue:

1. pulisce la cartella di output per far spazio ai nuovi file, usando la libreria `os`.
2. per ogni coppia di file in input contenendo l'analisi di allineamento eseguita usando i due tool:
 - (a) salva il nome del file json in ooutput seguendo l'analisi.
 - (b) ne compara le differenze.

2.1 librerie

- **re:** usato in **parsers.py** per dividere ogni linea tra: [id_sequenza, sequenza, posizioni] nelle funzioni di parsing.
- **hashlib:** usato in **utils.py** per creare l'hash da inserire come nome ai file JSON di output.
- **json:** usato in **utils.py** per elaborare ulteriormente l'output al fine di comparare i 2 tools di allineamento.
- **datetime:** usato in **utils.py** per prendere il tempo da inserire nell'hash, al fine di rendere unico l'output ed evitare sovrascritture.
- **os:** usato in **utils.py**, **main.py** per gestire input e output. Anche per pulire la cartella input nel main.

2.2 descrizione metodi

runClustal(inputFile, reference_id, nseq = 3): funzione che esegue il parsing del file di allineamento clustal (inputFile), esegue il parsing degli allineamenti e li salva nel file di output.

nseq serve al parser in quanto la classe **ClustalParser** richiede il numero di sequenze da elaborare nelle sue funzioni.

runMuscle(inputFile, reference_id, nseq = 3): funzione che esegue il parsing del file di allineamento muscle (inputFile), esegue il parsing degli allineamenti e li salva nel file di output.

nseq serve al parser in quanto la classe **ClustalParser** richiede il numero di sequenze da elaborare nelle sue funzioni.

parse(self, filename, reference=None, list=[]): legge il file in ingresso e restituisce il file in input suddiviso in: reference, sequenze, lunghezza sequenze.

save(alignment, analyzer, reference_id=None, tool=None, path=None): crea file json di output chiamato reference_id_hash-sha1 del tempo di inizio esecuzione concatenato con sequences_ids delle sequenze. .

3 Output

Come output vengono prodotti dei file json corrispondenti al nome reference_seq_hash.hexdigest(): l'id della reference accostata all'hash sha1 con encoding utf8 dell'id delle sequenze usate concatenate al tempo di inizio lavorazione (aggiunto per evitare di sovrascrivere l'output dello stessa serie di sequenze sia con muscle che con clustal).

l'oggetto JSON è così composto: { "reference": reference_Id,

"analyzed_sequences": insieme di sequenze analizzate,

"unmatches": altri oggetti JSON =

{ id (hash inizio mismatch + fine mismatch): "from": inizio mismatch,

"to": fine mismatch,

"sequences": lista sequenze in diverse tra loro in quell'intervallo }

}

4 Conclusioni