



Università degli Studi di Milano-Bicocca

Dipartimento di Informatica
Corso di Informatica - E3101Q
A.A. 2016-2019

Inferenza di alberi tumorali tramite Particle Swarm Optimization

Relazione della prova finale
Adrian David Castro Tenemaya, 816015
8 giugno 2019
Tutor: *Dott.* Ciccolella Simone
Relatore: *Prof.* Della Vedova Gianluca

Indice

1	Introduzione	7
1.1	Descrizione	7
1.2	Storia	7
1.3	Nozioni di biologia	7
1.3.1	La cellula	7
1.3.2	Il DNA	8
1.3.3	Cancro e tumore	9
1.4	Modello di sostituzione	9
2	Stato dell'arte	11
2.1	Introduzione	11
2.2	SiFit: inferring tumor trees from single-cell sequencing data under finite-sites models, 2017 [7]	11
2.2.1	Problema e soluzione	11

Premessa e ringraziamenti

Il presente lavoro è frutto del lavoro svolto come tirocinio all'interno dell'Università di Milano-Bicocca, e viene anche utilizzato come tesi finale ai fini del conseguimento della laurea in Informatica. È però necessario chiarire che il progetto in questione non sarà abbandonato nè una volta terminata la stesura di questa relazione, nè dopo il conseguimento della laurea. È mia intenzione contribuire al meglio delle mie possibilità in quello che ritengo essere uno dei campi con il quale mi sento più legato, sia a livello di interesse professionale, che a livello strettamente personale: la ricerca sul cancro. Secondo il *National Cancer Institute*, nel 2012 sono stati riportati 14.1 *milioni* di nuovi casi, e di questi, 8.2 *milioni* hanno portato alla morte [1]. I dati mostrano anche quelli che può sembrare all'apparenza una realtà discordante: il numero totale di morti per cancro è in crescita, ma il rapporto delle morti per individuo sta calando [3]. Nel 1990, 161 persone su 100.000 nel mondo sono morte a causa del cancro. Nel 2016, questo numero è calato a 134 su 100.000. Questo miglioramento è dovuto indubbiamente ad un numero molto elevato di fattori, tra cui l'aumento della qualità di vita ed un migliore sistema sanitario, ma è anche grazie alla crescita incessante della ricerca sul cancro, ed ai campi sui quali essa si appoggia. Lo sviluppo di algoritmi sempre più efficienti e performanti, e l'utilizzo di calcolatori super-veloci, ha permesso a questo settore di ricerca di ottenere dei considerevoli risultati.

Con questo progetto spero, quindi, di aver dato un contributo in questo settore, anche se in una percentuale minuscola.

Vorrei ringraziare mia mamma **Laura**, mio padre **José**, mia sorella **Valeria**, i miei fantastici nonni, e tutte le bellissime e meravigliose persone che hanno contribuito, in maniera diretta ed indiretta, a farmi appassionare all'informatica e, in questo caso, alla bioinformatica.

Prefazione

Il presente lavoro è stato svolto sotto la guida ed il supporto di AlgoLab, laboratorio presso il dipartimento di informatica dell’Università di Milano-Bicocca, che ha lo scopo di progettare, studiare, analizzare ed implementare algoritmi efficienti per problemi computazionali. Il tirocinio è cominciato il 22 Marzo 2019, ed è stato condotto per la maggior parte in maniera autonoma, da remoto. Il problema affrontato è l’*inferenza di progressioni tumorali* su dati single-cell, al fine di determinare l’ordine e la frequenza con cui le variazioni somatiche vengono acquisite durante una progressione tumorale. Spesso ciò è basato sulla “Infinite Sites Assumption”, dove le mutazioni possono solo essere acquisite, e mai perse. Lo stage si colloca nella ricerca del superamento di tale assunzione, utilizzando il modello della *filogenesi persistente*, dove ogni mutazione può essere persa al massimo una volta nell’intero albero. Più precisamente, si è investigata la tecnica *Particle Swarm Optimization*, un algoritmo di ottimizzazione di tipo euristico, ispirato al movimento degli sciami. I dati single-cell sono caratterizzati da un elevato tasso di errore e di valori mancanti: ciò rende inutilizzabili gli approcci noti in letteratura per i dati di *bulk sequencing*. In particolare, sono state analizzate quali strutture dati utilizzare per rendere l’algoritmo efficiente ed efficace, e quali operazioni considerare per inferire predizioni accurate.

Capitolo 1

Introduzione

1.1 Descrizione

Recenti sviluppi nel trattamento mirato di questo gruppo di malattie fa affidamento sull'accurata inferenza della progressione e dell'evoluzione del cancro [Ciccolella268243], rivelandosi una . Il cancro è la seconda causa più comune di morte [3], arrivando nel 2017 a contare il 17.08% delle morti nel mondo, per un totale di 8.93 *milioni* di decessi.

1.2 Storia

Era il 1869 quando venne isolato per la prima volta nella storia dell'umanità l'*Acido Desossiribonucleico*, anche conosciuto come *DNA*. Il pioniere di questa scoperta è Friedrich Miescher, medico e ricercatore nato in Svizzera nel 1844. Durante il processo di scoperta, Miescher aveva realizzato che nonostante avesse proprietà simili alle proteine, la nuova sostanza – il DNA – non lo era. Prima di isolare le cellule dal pus presente nelle bende chirurgiche dell'ospedale in cui lavorava, Miescher fu molto attento ad assicurarsi che il materiale che stava utilizzando fosse fresco e non contaminato. Fu solo più tardi, nel 1871, che il ricercatore iniziò a lavorare sullo sperma di salmone, una specie di pesce che affluiva numerosa durante il periodo autunnale nella città di Basel.

1.3 Nozioni di biologia

In questa sezione verranno introdotte nozioni inerenti al lavoro svolto, al fine di poter agevolarne la replicazione e la comprensione.

1.3.1 La cellula

Le cellule costituiscono le fondamenta di tutti gli organismi viventi. Il corpo umano è composto da trilioni di cellule. Esse danno forma al corpo, estraggono le sostanze nutritive dal cibo, convertono quelle sostanze nutritive in energia, ed hanno delle funzioni specifiche. Le cellule contengono anche il materiale ereditario del corpo, e possono fare copie di loro stesse [4]. Esse sono a loro volta costituite da diverse parti, tra le quali analizzeremo il nucleo e ciò che esso contiene, il DNA.

1.3.2 Il DNA

Il *DNA*, o *acido desossiribonucleico*, è il materiale ereditario negli umani e quasi tutti gli altri organismi viventi. Quasi ogni cellula presente all'interno del corpo umano ha lo stesso identico DNA. La maggior parte del DNA è situata all'interno del nucleo della cellula (dove è chiamato *DNA cellulare*), ma può trovarsi anche all'interno dei mitocondri, organi cellulari addetti alla respirazione cellulare. Le informazioni nel DNA sono conservate come un codice formato da quattro componenti chimici base (anche dette basi azotate): **adenina** (A) (Figura 1.2), **guanina** (G) (Figura 1.3), **citosina** (C) (Figura 1.4), e **timina** (T) (Figura 1.5). L'ordine, o la sequenza, di queste basi determina le informazioni disponibili per costruire e mantenere operativo un organismo.



Figura 1.1: Il DNA

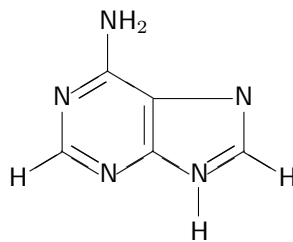


Figura 1.2: Adenina (A)

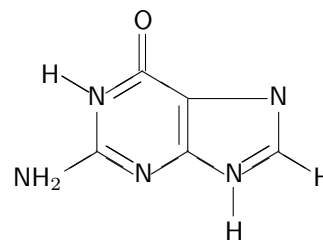


Figura 1.3: Guanina (G)

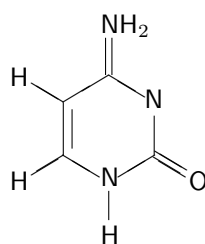


Figura 1.4: Citosina (C)

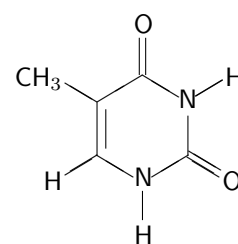


Figura 1.5: Timina (T)

Le basi del DNA si combinano tra di loro, A con T e C con G, in maniera tale da formare unità chiamate *coppie base*. Assieme ad una molecola di zucchero ed una di fosfato, le basi costituiscono quello che è definito un *nucleotide*. I nucleotidi sono disposti in due lunghi fili che formano una spirale chiamata *doppia elica*.

Un'importante proprietà del DNA è che si può replicare, o fare copie di se stesso. Uno qualunque dei due fili di DNA [5] può essere utilizzato nel processo di duplicazione per ottenere una copia identica del DNA di partenza. Questa è una fase cruciale nella divisione di una cellula, poiché la nuova copia di essa deve avere lo stesso identico DNA della cellula di origine.

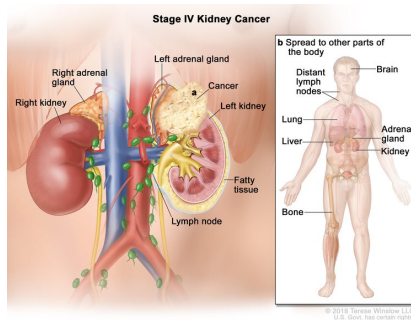


Figura 1.6: Cancro al rene

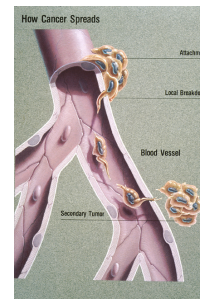


Figura 1.7: Diffusione del cancro

1.3.3 Cancro e tumore

La fase di divisione di una cellula è cruciale. Si stima che durante la replicazione, solo una base su 10^9 [2] sia errata. Vari fattori possono inoltre influenzare questa delicata fase, come l'esposizione ad agenti chimici ed irradiazione. Questi errori molto spesso sono corretti in vari modi, ma quando questo non basta, possono essere la causa scatenante che porta una cellula a diventare *cancerogena*. In generale, una cellula è cancerogena quando inizia a moltiplicarsi senza controllo. Quando questo processo avviene in un tessuto solido come un organo (Figura 1.6), muscolo od ossa, prende il nome di *tumore*. Ci sono due tipi di tumore: *maligno* (cancerogeno) e *benigno* (non cancerogeno). I primi possono invadere i tessuti circostanti nel corpo, e mentre crescono, alcune cellule possono viaggiare attraverso il sangue (Figura 1.7) o altri mezzi a formare delle *metastasi*, dei tumori secondari [6], mentre gli ultimi conservano le caratteristiche del tessuto di origine e non hanno la tendenza di invadere gli organi circostanti. Un tumore benigno non è quindi un cancro, ma solo una massa che può raggiungere dimensioni considerevoli, ma non si diffonde in altre parti del corpo.

1.4 Modello di sostituzione

In filogenetica il DNA può essere rappresentato come una sequenza di simboli, utilizzando le basi (sottosezione 1.3.2) corrispondenti alle posizioni degli allineamenti come caratteri.

AGTCCAGGACAT GGCATTCAATCA

Figura 1.8: Esempi di sequenze di DNA

La Figura 1.8 rappresenta un esempio di *modello di sostituzione*, un modello che in biologia descrive il processo per cui una sequenza di simboli cambia in un'altra, modificandone i tratti che rappresenta. In cladistica¹ viene utilizzato per rappresentare delle caratteristiche presenti, utilizzando il carattere "1", o assenti, utilizzando il carattere "0", in una specie.

10011 01110

Figura 1.9: Esempio di modello di sostituzione in cladistica

¹Metodo di classificazione dei viventi che si basa sul grado di parentela, ovvero sulla distanza nel tempo dell'ultimo progenitore comune

L'esempio in Figura 1.9 può ipoteticamente rappresentare due specie: la prima può digerire i latticini, non depone uova, è una creatura a sangue freddo, vola e sa nuotare; la seconda non può digerire i latticini, può deporre uova, è una creatura a sangue caldo, vola e non sa nuotare.

Lo stesso ragionamento può essere utilizzato per rappresentare le *mutazioni* acquisite o perse da una cellula nel corso della sua vita.

	<i>BBS4</i>	<i>CAMSAP1</i>	<i>DOCK3</i>	<i>EPHA10</i>	<i>EYA4</i>	<i>HIPK4</i>	<i>HIST1H2AG</i>	<i>INTS8</i>	<i>MAL2</i>	<i>MYOM3</i>	<i>OAZ3</i>	<i>PPIG</i>	<i>PTPRQ</i>	<i>RGS11</i>	<i>RYR3</i>	<i>SERPINF2</i>	<i>SMOC1</i>	<i>TTN</i>	<i>TUFT1</i>	<i>ZNF540</i>
1	1	1	0	0	0	1	0	1	1	1	1	1	1	1	0	0	1	0	0	
1	1	1	1	1	1	0	1	0	1	0	0	0	1	1	0	1	0	1	0	
0	0	0	0	0	0	1	0	1	1	1	1	1	0	1	1	0	1	0	1	
0	0	0	0	0	0	1	0	1	1	1	1	1	0	1	1	0	1	0	1	
0	0	0	0	0	0	1	0	1	1	1	1	1	0	1	1	0	1	0	1	
1	1	1	1	1	1	0	1	0	1	0	0	0	1	1	0	1	0	1	0	
0	0	0	0	0	0	1	0	1	1	1	1	1	0	1	1	0	1	0	1	
1	1	1	1	1	1	0	1	0	1	0	0	0	1	1	0	1	0	1	0	
0	0	0	0	0	0	1	0	1	1	1	1	1	0	1	1	0	1	0	1	

Figura 1.10: Esempio di modello di sostituzione con cellule (righe) e mutazioni (colonne)

Capitolo 2

Stato dell'arte

2.1 Introduzione

Con l'avvento delle tecnologie per il sequenziamento del DNA partendo da singole cellule (SCS), iniziano ad essere disponibili dati di alta qualità. Queste tecnologie forniscono il sequenziamento di dati da singole cellule, permettendo quindi di ricostruire l'albero filogenetico di una cellula. È però da tenere in considerazione l'alto tasso di errore associato a questo tipo di dati, innalzando di conseguenza il grado di difficoltà del processo di ricostruzione della filogenesi. In questo capitolo, analizzeremo le tecnologie già presenti che hanno affrontato questa sfida.

2.2 SiFit: inferring tumor trees from single-cell sequencing data under finite-sites models, 2017 [\[7\]](#)

2.2.1 Problema e soluzione

asdf

Bibliografia

- [1] *Cancer Statistics*. URL: <https://www.cancer.gov/about-cancer/understanding/statistics>.
- [2] Cooper GM. *The Cell: A Molecular Approach. 2nd edition - DNA Replication*. URL: <https://www.ncbi.nlm.nih.gov/books/NBK9940/>.
- [3] Max Roser e Hannah Ritchie. *OurWorldInData - Cancer*. Lug. 2015. URL: <https://ourworldindata.org/cancer>.
- [4] *What is a cell? - Genetics Home Reference - NIH*. URL: <https://ghr.nlm.nih.gov/primer/basics/cell>.
- [5] *What is DNA? - Genetics Home Reference - NIH*. URL: <https://ghr.nlm.nih.gov/primer/basics/dna>.
- [6] *What is the Difference Between Cancer and a Tumor? — Dana-Farber Cancer Institute*. URL: <https://blog.dana-farber.org/insight/2018/05/difference-cancer-tumor/>.
- [7] Hamim Zafar et al. «SiFit: inferring tumor trees from single-cell sequencing data under finite-sites models». In: *Genome biology* 18.1 (2017), p. 178. DOI: [10.1186/s13059-017-1311-2](https://doi.org/10.1186/s13059-017-1311-2).