



UNIVERSITÀ  
DEGLI STUDI  
FIRENZE

Scuola di Scienze Matematiche, Fisiche e Naturali  
Corso di Laurea in Informatica

APPLICAZIONI DELL'ALGORITMICA ALLA  
BIOLOGIA: EVOLUTIONARY TREES E  
CLUSTERING

APPLICATIONS OF ALGORITHMICS TO  
BIOLOGY: EVOLUTIONARY TREES AND  
CLUSTERING

MATTEO TORTOLI

MARIA CECILIA VERRI

Anno Accademico 2018-2019



---

## CONTENTS

---

1	Capitolo 1: La bioinformatica	7
1.1	Che cosa è la bioinformatica?	7
1.2	Aree di ricerca	9
1.2.1	Analisi dei genomi	9
1.2.2	Analisi di sequenze	9
1.2.3	Filogenetica	10
1.2.4	Bioinformatica strutturale	11
1.2.5	Espressione genica	12
1.2.6	Genetica delle popolazioni	12
1.2.7	Biologia dei sistemi	12
1.2.8	Data mining	12
1.2.9	Database biologici	12
1.2.10	Bioimmagini	12
1.3	Storia	12



---

## LIST OF TABLES

---



---

## LIST OF FIGURES

---





---

## CAPITOLO 1: LA BIOINFORMATICA

---

Per molti anni l'informatica è stata una scienza a sé stante, tuttavia negli ultimi decenni, grazie al progresso scientifico e tecnologico, sono nate nuove discipline chiamate genericamente **X-Informatics**. Queste sono il risultato dell'incontro tra l'informatica ed altre scienze di base (quali la biologia, la chimica, l'astronomia, la geologia etc) e tra queste citiamo la bioinformatica, la chemioinformatica, l'astroinformatica, la geoinformatica e così via. Anche se queste discipline sono diverse tra loro, ad esempio i dati raccolti in campo astronomico saranno di natura diversa rispetto quelli raccolti in campo biologico, condividono gli stessi obiettivi, come riportato nella pubblicazione [3] *X-Informatics: Practical Semantic Science*:

- Processamento ed estrazione delle informazioni
- Utilizzo trasparente ed efficiente dei dati in base al contesto scientifico, dalla raccolta, all'analisi fino alla catalogazione
- Integrazione di dati ottenuti tra sorgenti eterogenee
- Interazione con la raccolta dati adattata e personalizzata per l'utente
- Fornire supporto decisionale per l'utente, riducendo così i possibili errori e facilitando l'analisi dei risultati

Tra tutte queste discipline, risulta di particolare importanza la bioinformatica.

### 1.1 CHE COSA È LA BIOINFORMATICA?

Non esiste un unico modo con il quale definire la bioinformatica, infatti è possibile trovare definizioni diverse tra loro in quanto i professionisti non sempre concordano sulla portata del suo uso sia nel campo della biologia che dell'informatica. Tuttavia una possibile definizione è la seguente:

*La bioinformatica è un campo multidisciplinare della scienza che coinvolge la genetica, la biologia molecolare, l'informatica, la matematica e la statistica, rivolta a studiare sistemi biologici utilizzando metodi e modelli informatici e computazionali.*

Tra i vari obiettivi precedentemente elencati nell'introduzione al capitolo, va aggiunto quello che risulta l'obiettivo principale di questa disciplina, ovvero quello di aumentare la conoscenza di tutti quei processi di natura biologica.

A prescindere dalla natura del problema da affrontare, è possibile individuare un approccio standard, suddiviso in cinque steps:

- Studio ed analisi del problema da affrontare
- Collezionamento ed analisi di dati statistici a fronte di dati biologici in input
- Creazione di modelli ed uso di strumenti matematici che possano essere applicati al problema in esame, al fine di sviluppare un algoritmo
- Creazione, valutazione e test dell'algoritmo risolutivo del problema

Una parte fondamentale della bioinformatica consiste in esperimenti che generano dati ad alto throughput (high-throughput data), tra cui la misurazione dei modelli di espressione genica oppure la determinazione della sequenza genomica. Per **high-throughput data** si intendono quei dati biologici ottenuti tramite tecniche automatizzate e quindi non ottenibili attraverso metodi convenzionali. Il mining di questi dati può portare a nuove scoperte scientifiche non solo in campo biologico, ma anche medico, sia nel breve che nel lungo periodo.

Nel breve periodo, ad esempio grazie al *progetto genoma umano*<sup>1</sup>, si possono scoprire nuovi geni legati alle malattie e nuovi bersagli molecolari, ovvero quei processi biologici, intesi come proteine, recettori, pathway biochimici etc su cui si può intervenire per modificare il decorso di una malattia.

Nel lungo periodo sarà possibile scoprire eventuali reazioni avverse ai farmaci da individuo ad individuo in base a dei tests, al punto tale che,

<sup>1</sup> Il progetto genoma umano (Human Genome Project) è stato uno dei più grandi progetti scientifici degli ultimi anni. L'obiettivo era quello di ottenere la sequenza del genoma umano (e quindi il suo intero DNA) e identificare in esso i geni contenuti. Il progetto è cominciato nel 1990, per poi essere completato nel 2003 ed ulteriori ricerche sono ancora in corso.

grazie all'informazione genetica ottenuta attraverso strumenti informatici, sarà possibile personalizzare l'uso di farmaci, portando ad una migliore efficacia alla terapia individuale, riducendo o addirittura eliminando possibili effetti collaterali.

## 1.2 AREE DI RICERCA

Data la natura eterogenea dei dati biologici, la bioinformatica comprende un vasto numero di aree di ricerca in continua crescita, di seguito presentate.

### 1.2.1 *Analisi dei genomi*

Uno dei principali focus della bioinformatica riguarda l'analisi dei genomi<sup>2</sup> degli organismi il cui sequenziamento<sup>3</sup> è già stato completato, dal moscerino della frutta fino all'essere umano. L'analisi dei genomi è un'area di ricerca relativa non solo alla bioinformatica, ma anche alla genomica, ovvero quella disciplina che studia la struttura, il contenuto, la funzione e l'evoluzione del genoma. Ma perché analizzare i genomi? un gene viene sequenziato per conoscere la sua funzione ed eventualmente per modificarne la sua funzione. La conoscenza dell'intero genoma di un organismo fornisce le sequenze di tutti i suoi geni, permettendo così di identificare e manipolare i geni importanti che influenzano il metabolismo, la differenziazione e lo sviluppo cellulare e i processi patologici negli umani, animali e nelle piante.

L'obiettivo in questa area di ricerca per la bioinformatica è quello di identificare e alterare tutti quei geni che abbiano una particolare funzione biologica attraverso strumenti computazionali.

### 1.2.2 *Analisi di sequenze*

L'analisi delle sequenze di DNA, RNA o proteine è un processo mediante il quale tali macromolecole vengono sottoposte a dei metodi analitici al fine di capirne la struttura e le funzionalità.

<sup>2</sup> Per genoma si intende l'intero materiale genetico di un organismo, composto da DNA o RNA.

<sup>3</sup> Il sequenziamento è un processo mediante il quale viene determinata la struttura primaria delle macromolecole (composizione atomica e legami), quali il DNA, RNA e proteine.

Di particolare importanza risulta lo studio delle sequenze di DNA<sup>4</sup> che possono essere memorizzate in un computer attraverso una vasta varietà di metodi. La memorizzazione avviene attraverso l'uso di caratteristiche identificative di una determinata sequenza di DNA, ad esempio il nome di un gene oppure la fonte, dopodiché vengono salvate all'interno di database che prendono il nome di *database biologici* (vedere sottosezione 1.2.9).

L'analisi delle sequenze di DNA permette di venire a conoscenza dell'informazione genetica che viene trasportata all'interno di un suo segmento, grazie al quale, ad esempio, è possibile individuare i cambiamenti di un gene che possono causare una potenziale malattia.

Una volta estratto il DNA, vengono create migliaia e migliaia di copie di un singolo frammento, in quanto non è singolo frammento non risulterebbe sufficiente. Questi campioni vengono inseriti in dei macchinari chiamati *DNA Sequencer* che svolgono il compito di sequenziamento (del DNA) in modo automatico, infine i dati vengono raccolti ed analizzati.

Tra i vari algoritmi utilizzati per l'analisi delle sequenze risultano di particolare importanza gli **algoritmi di Clustering**, il cui obiettivo è quello di raggruppare i dati delle sequenze in modo veloce e preciso. I clusters giocano un ruolo chiave all'interno della bioinformatica, infatti data la grande quantità di dati da gestire e manipolare, è possibile raggrupparli in insiemi (clusters, appunto) secondo determinati criteri, di modo che gli elementi che sono simili tra di loro stiano nello stesso cluster mentre quelli che sono differenti risiedono in altri clusters. I principali algoritmi di Clustering nella bioinformatica sono il *Clustering gerarchico* e l'*algoritmo k-Means*, che verranno spiegati successivamente.

### 1.2.3 Filogenetica

La filogenetica è quel campo della bioinformatica e della biologia evolutiva<sup>5</sup> che studia le relazioni evolutive tra vari gruppi di organismi. Tradizionalmente si basava sul confronto delle caratteristiche morfologiche degli organismi, oggi invece si usano i dati molecolari delle sequenze, permettendo la costruzione di alberi filogenetici con molta accuratezza.

L'*albero evolutivo* o *albero filogenetico* è un diagramma che rappresenta le relazioni evolutive tra i vari organismi (spiegato successivamente), la cui

<sup>4</sup> Il sequenziamento del DNA consiste nel determinare la sequenza di nucleotidi all'interno di un suo frammento.

<sup>5</sup> La biologia evolutiva si occupa dello studio delle origini ed evoluzioni delle specie.

invenzione è dovuta a Charles Darwin, nel 1837. Una delle peculiarità è che si possono costruire alberi in base a dati genetici, genomici o morfologici, al fine di descrivere le relazioni che vi sono tra organismi viventi oppure tra specie estinte e specie viventi. Con questa nuova informazione è possibile dare una definizione alternativa alla filogenetica, ovvero come quella disciplina che si occupa di ricostruire ed analizzare gli alberi filogenetici.

Le applicazioni della filogenetica sono molteplici:

- **Bioinformatica e computing:** Gli alberi risultano una struttura dati di fondamentale importanza nel campo dell'informatica e degli algoritmi, infatti, molti di questi sviluppati per la filogenetica sono stati successivamente utilizzati anche in altri settori, in cui hanno trovato impiego.
- **Classificazione:** fornisce dei metodi di classificazione degli organismi viventi in modo efficiente ed accurato.
- **Forense:** può essere usata per valutare delle prove di DNA in casi giudiziari.
- **Identificazione dell'origine degli agenti patogeni:** insieme all'analisi delle sequenze (sottosezione 1.2.2) è possibile studiare ancora più a fondo gli agenti patogeni<sup>6</sup>, impedendo eventuali epidemie. Infatti, scoprire a quale specie vivente è collegato un determinato agente patogeno fornisce delle informazioni per scoprire quale potrebbe essere una eventuale forma di trasmissione.
- **Conservazione:** fornisce informazioni aggiuntive per impedire la scomparsa di animali o vegetali, garantendone, quindi, la loro conservazione.

#### 1.2.4 *Bioinformatica strutturale*

Le macromolecole (DNA, RNA e proteine) svolgono la loro funzione all'interno di un organismo grazie anche alla loro struttura nello spazio tridimensionale e quindi dell'avvolgersi delle loro sequenze di aminoacidi di cui sono composte, tuttavia conoscere tale struttura non è affatto banale, ad esempio se prendiamo in considerazione tutti e 20 gli aminoacidi delle proteine ed una proteina composta da 70 aminoacidi, otteniamo

---

<sup>6</sup> Gli agenti patogeni sono i virus, i batteri, etc...

$20^{70}(= 1.180591620717411303424 \times 10^{91})$  possibili strutture diverse che si possono ottenere (anche se la natura non ne ha selezionate così tante!). Ed è qui che entra in gioco la bioinformatica strutturale, ovvero come area di ricerca che si pone l'obiettivo di analizzare e ricostruire tramite algoritmi la struttura tridimensionale delle macromolecole.

La risorsa bioinformatica principale di questo settore è il *Proteine Data Bank*[12], un vero e proprio archivio di acidi nucleici e proteine in 3-D.

#### 1.2.5 *Espressione genica*

#### 1.2.6 *Genetica delle popolazioni*

#### 1.2.7 *Biologia dei sistemi*

#### 1.2.8 *Data mining*

#### 1.2.9 *Database biologici*

#### 1.2.10 *Bioimmagini*

### 1.3    STORIA

---

## BIBLIOGRAPHY

---

- [1] Allegretti M. (2014).  
*La biologia strutturale in movimento.*  
AIRInforma: AIRIcerca.  
<http://informa.airicerca.org/it/2014/10/04/biologia-strutturale-in-movimento/>
- [2] Bayat Ardeshir (2002).  
*Science, medicine, and the future: Bioinformatics.*  
BMJ, 324(7344), 1018–1022. DOI:  
10.1136/bmj.324.7344.1018
- [3] Borne, K. D.  
*X-Informatics: Practical Semantic Science.*  
<http://adsabs.harvard.edu/abs/2009AGUFMIN43E..01B>.  
George Mason University Fairfax VA, American Geophysical Union  
Fall Meeting, 12/2009. (Cited on page 7.)
- [4] Can T. (2013).  
*Introduction to Bioinformatics.*  
miRNomics: MicroRNA Biology and Computational Analysis.  
DOI:  
10.1007/978-1-62703-748-8\_4
- [5] Emery L.  
*Phylogenetics: An introduction.*  
European Bioinformatics Institute-European Molecular Biology  
Laboratory.  
<https://www.ebi.ac.uk/training/online/course/introduction-phylogenetics>  
Oxford.
- [6] Haque Sultan Omar (2016).  
*Phylogenetics.*  
Enciclopedia Britannica.  
<https://www.britannica.com/science/phylogenetics#accordion-article-history>

- [7] ICAR CNR: Istituto Di Calcolo E Reti Ad Alte Prestazioni.  
*Bioinformatica*.  
<https://www.icar.cnr.it/bio-informatica/>
- [8] Jiang, X., Lin, G., Liu, X., Zeng, X., Zou, Q.(2018).  
*Sequence clustering in bioinformatics: an empirical study*.  
Briefings in Bioinformatics.  
DOI:  
<https://doi.org/10.1093/bib/bby090>
- [9] Jones C. Neil and Pevzner A. Pavel.  
*An introduction to bioinformatics algorithms*.  
Massachusetts, Massachusetts Institute of Technology, 2004.
- [10] Mount W. David.  
*Bioinformatics: Sequence and Genome Analysis*.  
Cold Spring Harbor Laboratory Press, 2004.
- [11] National Human Genome Research Institute (2015).  
*DNA Sequencing Fact Sheet*.  
[https://www.genome.gov/about-genomics/fact-sheets/](https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Fact-Sheet)  
DNA-Sequencing-Fact-Sheet
- [12] Proteine Data Bank Website.  
*PDB*.  
<http://www.rcsb.org/> (Cited on page 12.)
- [13] Semple C., Steel M.  
*Phylogenetics*.  
Oxford, Oxford University Press, 2003.
- [14] Tramontano Anna (2003).  
*La grande scienza. Bioinformatica*.  
[http://www.treccani.it/enciclopedia/](http://www.treccani.it/enciclopedia/la-grande-scienza-bioinformatica_%28Storia-della-Scienza%29/)  
la-grande-scienza-bioinformatica\_%28Storia-della-Scienza%  
29/