



Universidad Nacional de San Agustín

Artificial Intelligence

Multiple Sequence Alignment using Particle Swarm Optimization

MSc. Vicente Machaca Arceda

2021

Content



Introduction

- Bioinformatics
- Problem
- Objective

Proposal

- Particle definition
- Movements
- Objective function
- Mutations

Experiments

- Datasets and params

Results

- PSO vs CLUSTALW

Conclusions

Overview



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

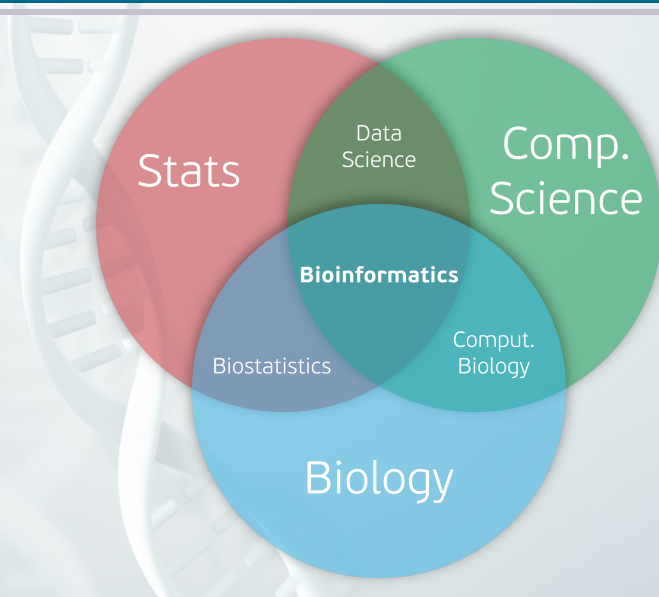
Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



Bioinformatics

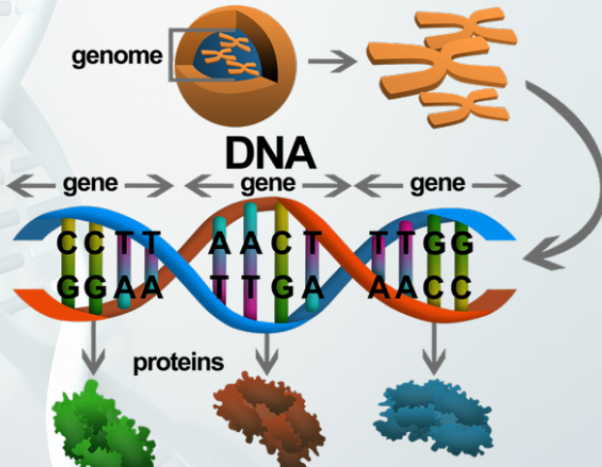
What is Bioinformatics?



According to Luscombe et al.: **Bioinformatics** involves the technology that uses computers for storage, retrieval, manipulation, and distribution of information related to biological macromolecules such as DNA, RNA, and proteins [1].



Bioinformatics is limited to sequence, structural, and functional analysis of genes and genomes and their corresponding products and is often considered **Computational molecular biology**. However, **Computational Biology** encompasses all biological areas that involve computation [2].





```
>J01859.1 Escherichia coli 16S ribosomal RNA, complete sequence
AAATTGAAGAGTTTGATCATGGCTCAGATTGAACGCTGGCGGCAGGCCTAACACATGCAAGTCGAACGGT
AACAGGAAGAAGCTTGCTCTTTGCTGACGAGTGGCGGACGGGTGAGTAATGTCTGGGAACTGCCTGATG
GAGGGGGATAACTACTGGAAACGGTAGCTAATACCGCATAACGTCGCAAGACCAAGAGGGGGACCTTCG
GGCCTCTTGCCATCGGATGTGCCAGATGGGATTAGCTAGTAGGTGGGGTAACGGCTCACCTAGGCGACG
ATCCCTAGCTGGTCTGAGAGGATGACCAGCCACACTGGAAC TGAGACACGGTCCAGACTCCTACGGGAGG
CAGCAGTGGGGAATATTGCACAATGGGCGCAAGCCTGATGCAGCCATGCCGCGTGTATGAAGAAGGCCTT
CGGGTTGTAAAGTACTTTTCAGCGGGGAGGAAGGGAGTAAAGTTAATACCTTTGCTCATTGACGTTACCCG
CAGAAGAAGCACCGGCTAACTCCGTGCCAGCAGCCGCGGTAATACGGAGGGTGCAAGCGTTAATCGGAAT
TACTGGGCGTAAAGCGCACGAGGCGGTTTGTTAAGTCAGATGTGAAATCCCCGGGCTCAACCTGGGAAC
TGCATCTGATACTGGCAAGCTTGAGTCTCGTAGAGGGGGGTAGAATCCAGGTGTAGCGGTGAAATGCGT
AGAGATCTGGAGGAATACCGGTGGCGAAGGCGGCCCCCTGGACGAAGACTGACGCTCAGGTGCGAAAGCG
TGGGGAGCAAACAGGATTAGATACCCTGGTAGTCCACGCGTAAACGATGTCGACTTGGAGGTTGTGCC
TTGAGGCGTGGCTTCCGGAGCTAACGCGTTAAGTCGACCGCTGGGGAGTACGGCCGCAAGGTTAAACT
CAAATGAATTGACGGGGGCCCGCACAAAGCGGTGGAGCATGTGGTTTAATTCGATGCAACGCGAAGAACCT
TACCTGGTCTTGACATCCACGGAAGTTTTT CAGAGATGAGAATGTGCCTTCGGGAACCGTGAGACAGGTGC
TGCATGGCTGTCGTGAGCTCGTGTTGTGAAATGTTGGGTTAAGTCCCGCAACGAGCGCAACCCCTTATCCT
TTGTTGCCAGCGTCCGGCCGGGAAC TCAAGGAGACTGCCAGTGATAAACTGGAGGAAGGTGGGGATGA
CGTCAAGTCATCATGGCCCCTTACGACCAGGGCTACACACGTGCTACAATGGCGCATACAAGAGAAGCGA
CCTCGCGAGAGCAAGCGGACCTCATAAAGTGCGTCGTAGTCCGATTGGAGTCTGCAACTCGACTCCATG
AAGTCGGAATCGCTAGTAATCGTGGATCAGAATGCCACGGTGAATACGTTCCCGGGCCTTGTACACACCG
CCCGTCACACCATGGGAGTGGGTTGCAAAAGAAGTAGGTAGCTTAACCTTCGGGAGGGCGCTTACCACCTT
TGTGATTCATGACTGGGGTGAAGTCGTAACAAGGTAACCGTAGGGGAACCTGCGGTTGGATCACCTCCTT
```

Figure: 16S ribosomal DNA of *Escherichia coli* with FASTA Format.

Overview



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions

Problem

Sequence alignment



No alignment

```
CGATGCTAGCGTATCGTAGTCTATCGTAC
      |      ||
ACGATGCTAGCGTTTCGTATCATCGTA
```

Aligned

```
-CGATGCTAGCGTATCGTAGTCTATCGTAC
||||| |||||
ACGATGCTAGCGTTTCGTA-TC-ATCGTA-
```

In the alignment process there could be substitutions, changes of residues and gaps. Gaps could cause by insertions or deletions.

Figure: No alignment versus alignment.

Problem

Sequence alignment



No gaps (10 matches)

```
a:  ATATTGCTACGTATATCAT
      |||||
b:  ATATATGCTACGTATCAT
```

With one gap (14 matches)

```
a:  ATAT-TGCTACGTATATCAT
      |||  |||||
b:  ATATATGCTACGTATCAT
```

With two gaps (16 matches)

```
a:  ATAT-TGCTACGTATATCAT
      |||  |||||  |||||
b:  ATATATGCTACG--TATCAT
```

Algorithms should take into account the possibility of introducing gaps. **Several alignments can be constructed** between two sequences.

Figure: Alignment and gaps.

Problem

Multiple Sequence alignment



```
RLA0_METVA  --MIDAKSEHKIAPWKIEEVNALKE LLKSANVIALIDMMEVPAVQLQEIRDK
RLA0_METJA   ---METKVKAHVAPWKIEEVKTLKGLIKSKPVVAIVDMMDVPAVQLQEIRDK
RLA0_PYRAB   -----MAHVAEWKKKEVEELANLIKSYPVIALVDVSSMPAYPLSQMRRL
RLA0_PYRHO   -----MAHVAEWKKKEVEELAKLIKSYPVIALVDVSSMPAYPLSQMRRL
RLA0_PYRFU   -----MAHVAEWKKKEVEELANLIKSYPVVALVDVSSMPAYPLSQMRRL
RLA0_PYRKO   -----MAHVAEWKKKEVEELANIIKSYPVIALVDVAGVPAYPLSKMRDK
RLA0_HALMA   MSAESERKTETIPEWKQEEVDAIVEMIESYESVGVVNIAGIPSRQLQDMRRD
RLA0_HALVO   MSESEVRQTEVIPQWKREEVDELVDFIESYESVGVVGAGIPSRQLQSMRRE
RLA0_HALSA   MSAAEQRTTEEVPEWKQEV AELVDLLETYDSVGVVNVGTGIPSKQLQDMRRG
RLA0_THEAC   -----MKEVSQQKKELVNEITORIKASRSVAIVDTAGIRTRQIQDIRGK
RLA0_THEVO   -----MRKINPKKKEIVSELAQDITKSKAVAIVDIKGVRTROMQDIRAK
RLA0_PICTO   -----MTEPAQWKIDFVKNLENEINSRKVAAIVSIKGLRNNEFQKIRNS
```

Figure: Example of Multiple Sequence Alignment (MSA) in amino acid sequences.

Overview



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



Propose Particle Swarn Optimization (PSO) to solve the Multiple Sequence Alignment (MSA) [3].

Overview



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



Leader:	W	G	K	V	-	-	N	V	D
	W	D	K	V	-	-	N	-	-
	S	-	K	V	G	G	N	-	-
Particle:	W	G	K	-	-	V	N	V	D
	-	W	D	K	-	-	-	V	N
	S	-	K	V	G	G	-	N	-

Leader:	5	6		
	5	6	8	9
	2	8	9	
Particle:	4	5		
	1	5	6	7
	2	7	9	

Figure: Example of particle representation.



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



$$distance = \frac{matchingGaps}{totalGaps} \quad (1)$$

$$crossPoint = rand(1, distance * length) \quad (2)$$

Leader:	5 6				
	5	6	8	9	
	2 8 9				
Particle:	4 5				
	1	5	6	7	
	2 7 9				

Particle	5				
	5	6	7		
	2	7	9		

Figure: Example of particle movement.



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions

ACGTCTGAT**A**CGCCGTAT**A**GTCTATCT
 | | | | | | | | | | | | | | |
----CTGAT**T**CGC---AT**C**GTCTATCT

Matches: $18 \times (+1)$

Mismatches: 2×0

Gaps: $7 \times (-1)$

Score = +11

Figure: Example of score in sequence alignment.

Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



The mutation operator inserts a gap in a random position in a random sequence inside a particle.

Introduction

Bioinformatics
Problem
Objective

Proposal

Particle definition
Movements
Objective function
Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



Table: Dataset used in the experiments.

<i>Dataset</i>	min. length	max. length	num. bases
S6	8	17801	153
S7	457	457	8
S8	7	10	5



1	A	T	G	C	A	A	G			
2	T	A	A	G	T	C	A	A	G	T
3	A	T	G	C	A	A	C	T		
4	T	A	A	G	T	C	A	T	A	
5	A	T	G	G	A	T	T	C		

Figure: Sequences of S8 dataset.

Table: Params used in the experiments.

Param	Value
Iterations	30
Num. of particles	25
Mutation probability	0.2
Gaps	30%
Num. of experiments test	10



Introduction

Bioinformatics

Problem

Objective

Proposal

Particle definition

Movements

Objective function

Mutations

Experiments

Datasets and params

Results

PSO vs CLUSTALW

Conclusions



Table: Score comparison of PSO and CLUSTALW.

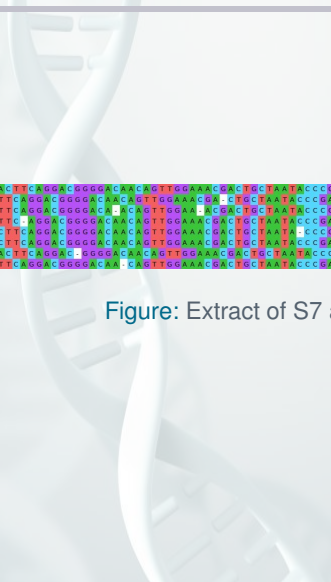
<i>Dataset</i>	PSO-mutation	PSO	CLUSTALW
S6	12678	10012	18045
S7	11105	9054	12564
S8	32	28	49



1	-	A	T	G	-	C	-	A	A	G	
2	T	A	A	G	T	C	A	A	G	T	
3	-	A	-	T	G	C	A	A	C	T	
4	T	A	A	G	T	C	A	T	-	A	
5	-	A	T	G	G	A	T	T	C	-	

1	A	T	G	C	A	A	G	-	-	-	
2	-	T	A	A	G	T	C	A	A	G	T
3	A	T	G	C	A	A	C	T	-	-	-
4	-	T	A	A	G	T	C	A	T	A	-
5	A	T	G	G	A	T	T	C	-	-	-

Figure: Left: S8 alignment with PSO-mutation. Right: S8 alignment with CLUSTALW.



```
1 -ACCTCAGGACGGGGACAACAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGGAGAAAGAGCTTGCCTCTGATTAGCTAGT
2 ACTTCAGGACGGGGACACAGTTGGAAACGA-CTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
3 ACTTCAGGACGGGGACACAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
4 GTTCAGGACGGGGACACAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
5 ACTTCAGGACGGGGACACAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
6 ACTTCAGGACGGGGACACAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
7 -ACTTCAGGAC-GGGGACACAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
8 ACTTCAGGACGGGGACAA-CAGTTGGAAACGACTGCTAATACCCGATGTGCCGCAAGGTGAAACCCTAAATTGGCCCTGAAGAAAGAGCTTGCCTCTGATTAGCTAGT
```

Figure: Extract of S7 alignment using PSO-mutation

Conclusions



This thesis proposed the use of PSO to solve the MSA problem. The author used a set of datasets from NCBI.

The author proposed a mutation operator to avoid local solutions. This operator just inserts a gap.

The score of PSO was acceptable and very similar to CLUSTALW. Currently, there is more research on this topic.



- [1] N. M. Luscombe, D. Greenbaum, and M. Gerstein, “What is bioinformatics? a proposed definition and overview of the field,” *Methods of information in medicine*, vol. 40, no. 04, pp. 346–358, 2001.
- [2] J. Xiong, *Essential bioinformatics*. Cambridge University Press, 2006.
- [3] F. B. R. Zablocki *et al.*, “Multiple sequence alignment using particle swarm optimization,” Ph.D. dissertation, University of Pretoria, 2009.

