# Estimation of Distribution Algorithms for Combinatorial Problems in Bioinformatics

#### Roberto Santana

Intelligent Systems Group
The University of the Basque Country
San Sebastian, Spain
roberto.santana@ehu.eus

(Lyon, October 2, 2025)





### Index

- TagSNP selection problem
- Estimation of distribution algorithms (EDAs)
- Transfer learning in EDAs
- EDAs for TagSNP selection

### TagSNPs Problem

#### Single Nucleotide Polymorphisms

More than 99.9% sequence of DNA is similar between different populations

#### SNiPs or SNPs =

sites of variation in the genome (spelling mistakes)

```
Karen AGCTTGAC TCC ATGATGATT
Debo AGCTTGAC GCCATGATGATT
Jose AGCTTGAC TCC TGATGATT
Thomas AGCTTGAC GCCC TGATGATT
Anupriya AGCTTGAC TCC ATGATGATT
Michelle AGCTTGAC TCC TGATGATT
Michelle AGCTTGAC GCCC TGATGATT
AGCTTGAC GCCC TGATGATT
AGCTTGAC GCCC TGATGATT
AGCTTGAC GCCC TGATGATT
```

Most of the variations come given by single nucleotide polymorphism (SNPs)



### TagSNPs Selection Problem

- If an individual carries out a SNP then, with high probability will carry out another "close" SNP
- The value of one SNP can be used to predict the value of other SNP
- TagSNP selection problem:
   Given a set of SNPs, find the minimum subset that can be used to predict the rest of SNPs
- Prediction ability is given in terms of correlation

### TagSNPs Selection Problem

- Prediction ability is given in terms of correlation threshold  $r_{min}$
- Single Marker: A SNP  $s_i$  tags another SNP  $s_j$  if  $r_{i,j} > r_{min}$
- Multiple Marker: A subset of SNPs T tags a single SNP  $s_i$ , if  $r_{T,i} > r_{min}$

#### We consider marker subsets with $|T| \le 2$

Santana, R., Mendiburu, A., Zaitlen, N., Eskin, E., and Lozano, J. A. (2010). Multi-marker tagging single nucleotide polymorphism selection using estimation of distribution algorithms. Artificial Intelligence in Medicine, 50(3), 193-201.



### TagSNPs Selection Problem

#### Codification

• 
$$\mathbf{x} = (x_1, \dots, x_n)$$
 where

$$x_i = \begin{cases} 1 & \text{if } s_i \text{ is a tagSNP} \\ 0 & \text{otherwise} \end{cases}$$

#### Optimization problem

$$\min f(\mathbf{x}) = \sum_{i=1}^{n} x_i$$

subject to x tags all the SNPs

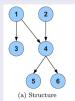
### **Datasets**

Index	Population	$\mid n.ind.$	SNPs	n'	pairs	triples
1	Adygei	15	1388	1048	3445	203640
<b>2</b>	Balochi	24	1399	1020	3486	220892
3	BantuKenya	11	1338	800	2183	89372
4	BantuSouthAfrica	8	1379	905	2553	127780
5	Basque	24	1362	1042	3956	252296
6	Bedouin	45	1433	1015	3147	182836
7	Bengali	15	1390	1007	3374	211439
8	BiakaPygmy	23	1322	700	2071	78136
9	Brahui	24	1368	975	2770	140452
10	Burusho	23	1388	1016	3272	197048
11	CEU	60	1386	1018	3227	191730
12	Cambodian	8	1335	1067	3476	219632
13	Colombian	7	1175	1008	4467	264906
14	Dai	10	1267	982	3762	212758
15	Daur	10	1248	956	3372	187016
16	Druze	41	1377	1019	3306	190120
17	French	28	1418	1068	3596	229216
18	Han	34	1273	946	3527	211194
19	Han-NChina	10	1230	954	3270	188553
20	Hazara	22	1329	968	2868	159534

Probabilistic graphical models

# Probabilistic graphical model

#### Graphical models



$X_1$	1	1	1	2	2	2
$X_2$	1	2	3	1	2	3
$P(X_4=1 \mid X_1,X_2)$	$\theta_{411}$	$\theta_{421}$	$\theta_{431}$	$\theta_{441}$	$\theta_{451}$	$\theta_{461}$
$P(X_4=2 \mid X_1,X_2)$	$\theta_{412}$	$\theta_{^{422}}$	$\theta_{432}$	$\theta_{442}$	$\theta_{\scriptscriptstyle 452}$	$\theta_{\scriptscriptstyle 462}$
$P(X_4=3 \mid X_1,X_2)$	$\theta_{413}$	$\theta_{423}$	$\theta_{433}$	$\theta_{443}$	$\theta_{453}$	$\theta_{463}$
$P(X_4=4 \mid X_1,X_2)$	$\theta_{414}$	$\theta_{424}$	$\theta_{434}$	$\theta_{444}$	$\theta_{\scriptscriptstyle 454}$	$\theta_{464}$
$P(X_4=5 \mid X_1,X_2)$	$\theta_{415}$	$\theta_{425}$	$\theta_{435}$	$\theta_{445}$	$\theta_{455}$	$\theta_{465}$

(b) Conditional probability table

A <u>probabilistic graphical model</u> for  $\mathbf{X} = (X_1, X_2, \dots, X_n)$  encodes a graphical factorization of a joint probability distribution  $p(\mathbf{x})$ 

- It has two components:
  - ullet A structure  $\mathcal S$  (e.g. directed acyclic graph for Bayesian networks).
  - A set of local marginal probability values.
- S represents a set of conditional independence assertions between the variables.



# Optimization approaches

#### Optimization problem

$$\hat{f} = max_{x_i \in \{1, \dots, K\}^n} f(\mathbf{x})$$

#### **EDAs**

- Population based evolutionary algorithms
- Use probabilistic models instead of genetic operators
- At each generation a probabilistic model of the selected population is learnt
- The probabilistic model is used to sample new solutions





Graphical models in optimization

### Probabilistic distributions in EDAs

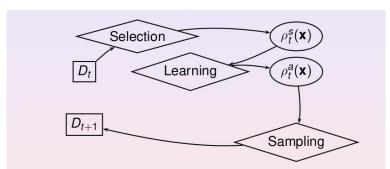


Figure: Joint probability distributions determined by the components of an EDA.  $D_t$ ,  $D_{t+1}$ : populations at generation t and t+1;  $\rho_t^s(\mathbf{x})$ ,  $\rho_t^a(\mathbf{x})$ : Joint probability distributions determined by selection and the probabilistic model approximation.





### Pseudocode of UMDA

#### UMDA (Muhlenbein and Paas:1996)

- $\bigcirc$   $D_0 \leftarrow$  Generate M solutions randomly
- 2 / = 1
- **3** do {
- $D_{l-1}^s \leftarrow \text{Select } N \leq M \text{ solutions from } D_{l-1}$  according to a selection method
- $p_l(\mathbf{x}) = \prod_{i=1}^n p(x_i|D_{l-1}^s) \leftarrow \text{Estimate the joint}$  probability of the selected solutions
- $D_l \leftarrow \text{Sample } M \text{ solutions (the new population)}$ from  $p_l(\mathbf{x})$
- } until A stop criterion is met



Graphical models in optimization

# EDAs: other probabilistic model used

#### Markov-chain model

$$p_{MK}(\mathbf{x}) = p(x_1, \dots, x_{k+1}) \prod_{i=k+2}^{n} p(x_i \mid x_{i-1}, \dots, x_{i-k})$$

#### Tree model

$$p_{Tree}(\mathbf{x}) = \prod_{i=1}^{n} p(x_i \mid pa(x_i))$$

#### Mixture of trees model

$$p_{MT}(\mathbf{x}) = \sum_{i=1}^{m} \lambda_{j} p_{Tree}^{j}(\mathbf{x})$$

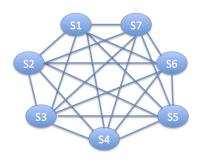


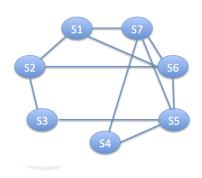
### EDA approaches: Tree model

#### Classical and restricted tree-EDA approaches

- Tree-EDA: The tree is learnt from the matrix of mutual information learned from data
- Tree-EDA<sup>r</sup>: The tree is learnt from the data constrained to known interaction of problem instance

### EDA for tagSNP selection problem



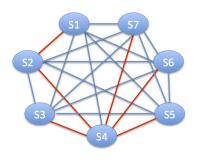


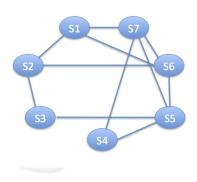
TreeEDA

TreeEDA<sup>r</sup>



### EDA for tagSNP selection problem



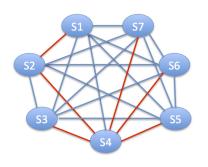


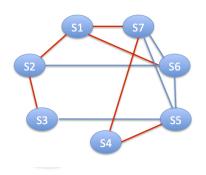
TreeEDA

TreeEDA<sup>r</sup>



### EDA for tagSNP selection problem





TreeEDA

TreeEDA<sup>r</sup>



#### **Datasets**

#### Characteristics of the datasets

No. of datasets: 58

No. of SNPs: 780–1089

No. of pairs: 1937–4467

No. of triples: 76007–264906

#### Algorithm Parameters

Pop. Size: 5000

Selection: 15

Repetitions: 30

Generations: 100



# Preliminary Results

I	Best			mean			
	UMDA	Tree	TreeR	UMDA	Tree	TreeR	
1	362	285	285	367.63	291.30	288.13	
2	365	302	300	371.17	304.50	<u>303.73</u>	
3	344	284	283	349.73	285.97	285.53	
4	372	305	306	378.47	308.67	308.33	
5	354	283	279	360.80	286.53	284.40	
6	376	299	298	381.63	302.37	301.73	
7	341	276	273	346.93	281.27	280.00	
8	325	265	264	329.23	265.70	265.43	
9	364	296	294	374.53	299.37	298.40	
10	351	288	288	359.30	291.43	291.50	
11	360	293	294	370.00	299.43	297.50	
12	341	270	270	346.93	276.00	273.97	
13	286	235	233	294.20	238.90	<u>236.30</u>	
14	325	263	261	332.37	266.90	263.13	
15	325	264	264	332.90	267.63	267.67	
16	361	291	291	369.90	297.37	296.50	
17	360	284	284	374.17	290.70	289.07	

#### Results

#### TreeEDA vs TreeEDA<sup>r</sup>

- Average results: TreeEDA<sup>r</sup> obtains better results in 43 from 58 datasets
- Best results: TreeEDA<sup>r</sup> obtains better results in 26 from 58 datasets and the same results in 15 of them

### Solving new instances

Instance 1

Instance 2

Instance 3

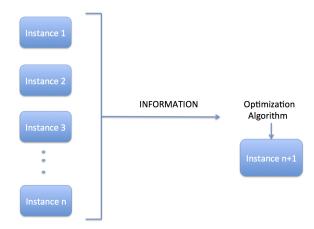
.

Instance r

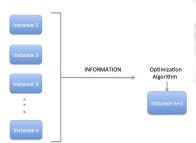
Instance n+1



### Solving new instances



# Transfer Learning

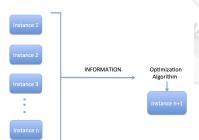


#### Questions

- What information extract from the problems?
- How to define relatedness between problems?
- How to transfer the available information?



# Transfer Learning



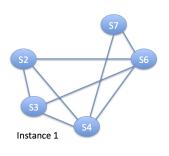
#### Transfer Learning for Optimization

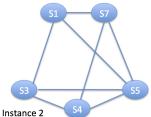
- Transfer of solutions
- Structural transfer
- Behavioral or algorithmic transfer
- Combined transfer

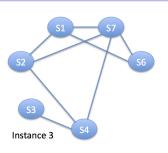
### Transfer Learning in the TagSNP problem

#### Structural Transfer

- The possible dependencies between the variables are transfered
- Four different transfer strategies are evaluated





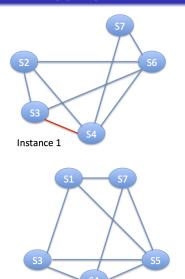


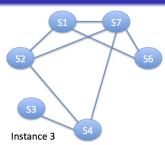




Transfer Learning for TagSNPs with EDAs

Instance 2

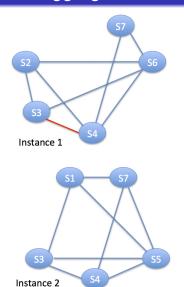


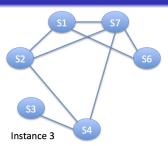






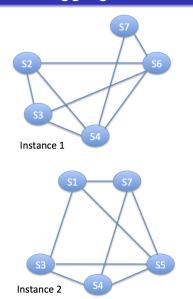


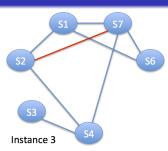


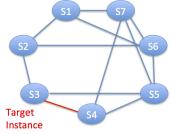




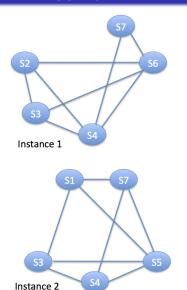


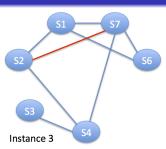


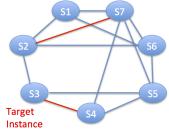




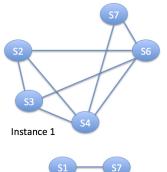


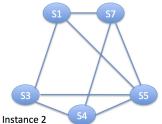


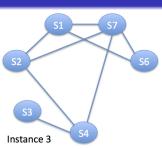








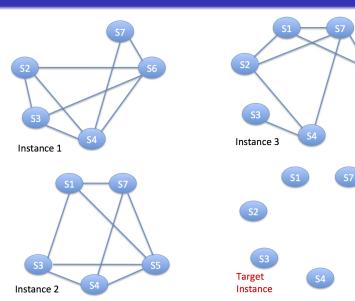






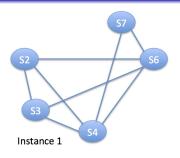


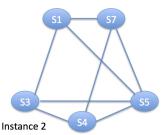
### Tn2: Pure Transfer

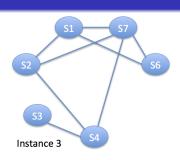


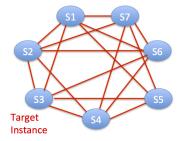
Transfer Learning for TagSNPs with EDAs

### Tn2: Pure Transfer



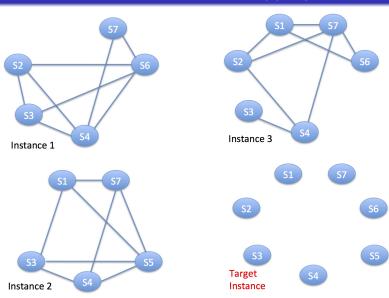




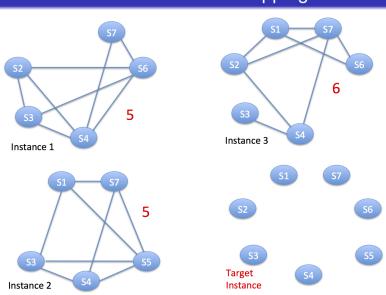




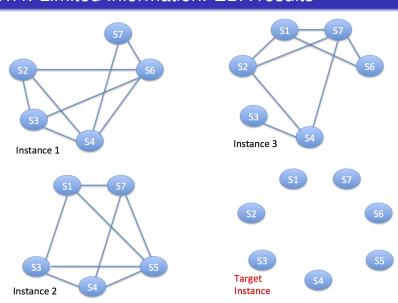
### Tn3: Limited Information. Overlapping Variables



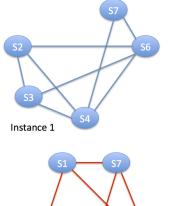
# Tn3: Limited Information. Overlapping Variables

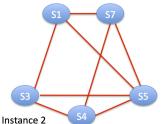


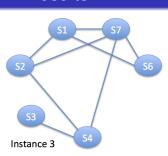
### Tn4: Limited Information. EDA results

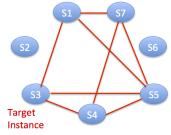


### Tn4: Limited Information. EDA results



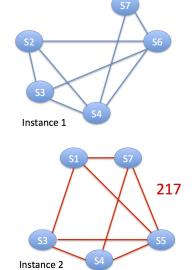




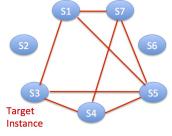




## Tn4: Limited Information. EDA results



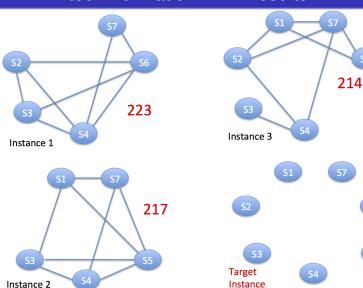






### Transfer Learning for TagSNPs with EDAs

## Tn4: Limited Information. EDA results



# Transfer Learning in the tagSNP Problem. Tn1 and Tn2

#### Tn1

 A matrix of interactions for problem i is created by combining the interactions between any pair of SNPs in problem i that were detected in any of the other SNP datasets

#### Tn2

 Identical to Tn1 but information about the interactions in problem i is not included in the aggregation matrix.



## Transfer Learning in tagSNP Problem Tn3 and Tn4

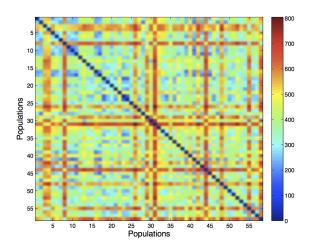
#### Tn3

- The most related problems to problem i are selected as those that share the largest set of common variables (SNPs) with it.
- The aggregation matrix is formed using the matrix of interactions from this subset of problems.

#### Tn4

- The interaction matrix of each problem j are used to solve problem i
- The best interaction matrices for problem i are aggregated

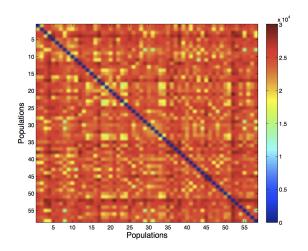
## Similarities between problems (share SNPs)







# Similarities between problems (structural similarity)







## Transfer Learning Results

mean			I	Best				mean			
UMDA	Tree	TreeR		Tn1	Tn2	Tn3	Tn4	Tn1	Tn2	Tn3	Tn4
367.63	291.30	288.13	1	281	281	280	283	285.60	285.20	287.00	287.33
371.17	304.50	303.73	2	296	298	299	297	300.13	300.97	302.93	302.00
349.73	285.97	285.53	3	281	282	283	284	284.27	284.60	285.80	285.13
378.47	308.67	308.33	4	304	304	306	306	306.47	306.17	308.60	309.03
360.80	286.53	284.40	5	276	278	282	282	281.53	281.47	286.93	286.00
381.63	302.37	301.73	6	295	297	295	294	299.23	298.97	300.50	299.00
346.93	281.27	280.00	7	271	269	271	271	275.60	275.00	277.97	275.80
329.23	265.70	265.43	8	263	263	265	264	265.07	264.63	266.43	265.73
374.53	299.37	298.40	9	291	293	291	292	294.80	295.73	295.47	296.17
359.30	291.43	291.50	10	284	285	283	284	287.33	288.20	287.80	287.80
370.00	299.43	297.50	11	291	291	289	288	296.50	295.67	295.33	294.23
346.93	276.00	273.97	12	265	266	270	269	271.07	271.70	274.67	274.20
294.20	238.90	236.30	13	230	229	229	232	233.57	233.47	235.83	236.80
332.37	266.90	263.13	14	257	255	260	258	261.03	261.27	263.37	262.83
332.90	267.63	267.67	15	260	261	266	261	264.73	264.47	268.17	264.97
369.90	297.37	296.50	16	290	286	289	288	293.63	292.17	292.40	293.60
374.17	290.70	289.07	17	281	279	282	284	285.80	286.10	287.27	288.63
331.87	262.37	261.50	18	255	255	255	256	258.13	258.43	258.63	259.17
336.60	269.43	268.00	19	261	262	260	262	265.90	266.30	266.03	267.00
366.77	296.37	294.00	20	287	288	289	291	292.40	292.03	293.33	295.03
332.33	263.27	260.70	21	253	254	254	253	257.27	257.27	258.77	257.67
340.97	278.23	277.27	22	270	270	272	271	274.53	274.67	275.80	274.47
327.50	261.23	260.00	23	253	254	255	254	257.10	256.43	257.93	257.00
360.00	282.43	280.80	24	271	272	274	272	277.37	277.20	278.97	277.33





# Transfer Learning Results

- Tn1 improves all the results of TreeEDA and TreeEDA<sup>r</sup>
- Tn2 improves in all but one
- Tn3 improves in 91% and 79% TreeEDA and TreeEDA<sup>r</sup> respectively
- Tn4 improves in 90% and 76% TreeEDA and TreeEDA<sup>r</sup> respectively

### Conclusions

- Transfer learning is an attractive alternative for EAs and EDAs
- Simple transfer learning can dramatically improve the results of optimization algorithms
- It is a new area in optimization that deserves much attention

# Estimation of Distribution Algorithms for Combinatorial Problems in Bioinformatics

#### Roberto Santana

Intelligent Systems Group
The University of the Basque Country
San Sebastian, Spain
roberto.santana@ehu.eus

(Lyon, October 2, 2025)





## Protein problems addressed

#### Protein folding in HP model

To find the HP configuration with lowest energy

#### Side chain placement problem

 To find the optimal positioning of the side chain with respect to a given backbone

#### Protein design

 To find the sequence that has the lowest energy with respect to a protein structure



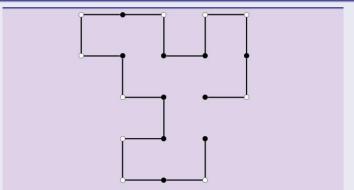


Transfer Learning for TagSNPs with EDAs

EDAs for protein problems

## HP and functional model protein

# Example





Transfer Learning for TagSNPs with EDAs

EDAs for protein problems

## Results in the two-dimensional lattice

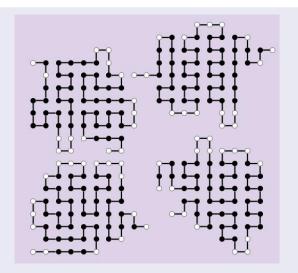




Figure: Ontimal solution and three sub-ontimal solutions for the s7