



## **Probability and Random Processes**

### **Project Report on**

### **“Expanding From Discrete To Continuous Estimation Of Distribution Algorithms”**

#### **GROUP - MEMBERS**

- |                                  |                   |
|----------------------------------|-------------------|
| <b>1. Shriyanshi Srivastava</b>  | <b>9919103158</b> |
| <b>2. Prakash Singh Rautela</b>  | <b>9919103160</b> |
| <b>3. Devendra Singh Chauhan</b> | <b>9919103161</b> |
| <b>4. Sachin Kumar Chauhan</b>   | <b>9919103175</b> |

**SUBMITTED TO:**     ***DR. AMITA BHAGAT***

## **Acknowledgment**

**We would like to start this project by expressing our special thanks of gratitude to our mentor, *Dr. Amita Bhagat* for giving us the opportunity to study and present our findings on the topic “Discrete and Continuous Probability Distribution”.**

**This project helped us to understand different areas of real life application of probability distribution. We would like to extend our sincere and heartfelt gratitude towards our mentor who has helped us in this endeavor with her active guidance, cooperation and encouragement.**

**At last, we would like to thank all our friends and family who directly or indirectly helped us in completing this project report.**

# **Index**

**1. Introduction**

**2. Objective**

**3. The Idea framework**

**4. Probability density structure search algorithms**

**5. Probability density functions**

**6. Conclusion**

## Introduction

Estimation of Distribution Algorithms (EDAs) is a new phase of Evolutionary Computation. Algorithms in evolutionary optimization guide their search through statistics based on a vector of samples, often called a population. By using this stochastic information, non-deterministic induction is performed in order to attempt to use the structure of the search space and thereby aid the search for the optimal solution. In order to perform induction, these samples are combined so as to generate new solutions that will hopefully be closer to the optimum. As this process is iterated, convergence is intended to lead the algorithm to a final solution.

In the genetic algorithm [11, 14] and many variants thereof, values for problem variables are often exchanged and subsequently individually adapted. Another way of combining the samples is to regard them as being representative of some probability distribution. Estimating this probability distribution and sampling more solutions from it, is a global statistical type of inductive iterated search. Such algorithms have been proposed for discrete spaces [2–4, 12, 13, 15, 17, 19, 21], as well as in a limited way for continuous spaces [5, 10, 15, 22, 23]. An overview of this field has been given by Pelikan, Goldberg and Lobo [20].

## Objective

Our goal in this paper is to apply the search for good probability density models to continuous spaces. To this end, we formalize the notion of building and using probabilistic models in a framework named IDEA. We show how we can adjust existing techniques to be used in the continuous case. We thereby define evolutionary optimization algorithms. Using a set of test functions, we validate their performance.

The remainder of this paper is organized as follows:

In section 3, we present the IDEA framework. In section 4, we describe a few existing algorithms that build and use probabilistic models. In section 5, we state some derivations of probability density functions (pdfs). Our final conclusions are drawn in section 6.

## The IDEA Framework

We write  $a = (a_0, a_1, \dots, a_{|a|-1})$  for a vector  $a$  of length  $|a|$ . The ordering of the elements in a vector is relevant. We assume to have  $l$  random variables available, meaning that each sample point is an  $l$  dimensional vector. We introduce the notation  $a_c = (a_{c0}, a_{c1}, \dots, a_{c|c|-1})$ . Let  $L = (0, 1, \dots, l-1)$  be a vector of  $l$  numbers and let  $Z = (Z_0, Z_1, \dots, Z_{l-1})$  be a vector of  $l$  random variables. We assume that we have an  $l$  dimensional cost function  $C(z_L)$  which without loss of generality we seek

to minimize. Without any prior information on  $C(z_L)$ , we might as well assume a uniform distribution over  $Z$ . Now denote a probability distribution that is uniformly distributed over all  $z_L$  with  $C(z_L) \leq \theta$  and that has a probability of 0 otherwise, by  $P_\theta(Z)$ . In the discrete case we have:

$$P^\theta(\mathcal{Z})(z(\mathcal{L})) = \begin{cases} \frac{1}{|\{z'(\mathcal{L}) | C(z'(\mathcal{L})) \leq \theta\}|} & \text{if } C(z(\mathcal{L})) \leq \theta \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Note that if we find  $P_{\theta^*}(Z)$  where  $\theta^* = \min_{z_L} \{C(z_L)\}$ , a single sample drawn from  $P_{\theta^*}(Z)$  provides an optimal solution  $z_L^*$ . A probability distribution is made up of a probability density structure (pds) and a probability density function (pdf) for each element in the pds. In graphical models literature, a pds is also called a factorization. Let  $ab$  be the splicing of  $a$  and  $b$  such that the elements of  $b$  are placed behind the elements of  $a$ , giving  $|ab| = |a| + |b|$ . Using graphical modelling [9], we can denote any non-clustered pds with conditional probabilities  $P(Z_a|Z_b) = P(Z_{ab})/P(Z_b)$ . We let  $\pi(\cdot)$  be a function that returns a vector  $\pi(i) = (\pi(i)_0, \pi(i)_1, \dots, \pi(i)|\pi(i)|-1)$  of indices denoting the variables that  $Z_i$  is conditionally dependent on. We call the graph that results when taking the  $Z_i$  as nodes and having an arc from node  $Z_i$  to node  $Z_j$  if and only if  $i \in \pi(j)$ , the pds graph. The only required condition to be able to express any non-clustered pds using conditional probabilities, is that this graph needs to be acyclic. By using a permutation vector  $\omega$ , the definition of a pds  $(\pi, \omega)$  that models conditional factorizations, can be formalized as follows:

$$P_{(\pi, \omega)}(\mathcal{Z}) = \prod_{i=0}^{l-1} P(Z_{\omega_i} | Z_{\langle \pi(\omega_i) \rangle}) \quad (2)$$

such that  $\forall i \in \mathcal{L} \langle \omega_i \in \mathcal{L} \wedge \forall k \in \mathcal{L} - (i) \langle \omega_i \neq \omega_k \rangle \rangle$   
 $\forall i \in \mathcal{L} \langle \forall k \in \pi(\omega_i) \langle k \in \{\omega_{i+1}, \omega_{i+2}, \dots, \omega_{l-1}\} \rangle \rangle$

Each  $P(Z_{\omega_i} | Z_{\pi(\omega_i)})$  from equation 2 is a special case multivariate conditional pdf. This means that any such conditional pdf along with a pds  $(\pi, \omega)$  defines a probability distribution over  $Z$ . In general, we denote a pds by  $f$ . The pds is constrained to be of a certain form. For instance, in the case of equation 2, the constraints impose  $f$  to describe a directed acyclic graph. We denote the constrained space of all possible structures by  $C$ . A probability distribution over  $Z$  is then formally denoted by  $P_f(Z)$ ,  $f \in C$ . Denote the largest function value of a selection of samples at iteration  $t$  by  $\theta_t$ . We find a pds and estimate each pdf to best approximate  $P_{\theta_t}(Z)$ . We can then sample from the resulting probability distribution to get more

samples. By formalizing this rationale in an iterative algorithm, we define the Iterated Density Estimation Evolutionary Algorithm (IDEA):

IDEA( $n, \tau, m, \text{sel}(), \text{rep}(), \text{ter}(), \text{sea}(), \text{est}(), \text{sam}()$ )	
Initialize an empty vector of samples Add and evaluate $n$ random samples	$\mathcal{P} \leftarrow ()$ <b>for</b> $i \leftarrow 0$ <b>to</b> $n - 1$ <b>do</b> $\mathcal{P} \leftarrow \mathcal{P} \sqcup \text{NewRandomVector}()$ $c[\mathcal{P}_i] \leftarrow C(\mathcal{P}_i)$
Initialize the iteration counter Iterate until termination Select $\lfloor \tau n \rfloor$ samples Set $\theta_t$ to the worst selected cost	$t \leftarrow 0$ <b>while</b> $\neg \text{ter}()$ <b>do</b> $(z^0 \langle \mathcal{L} \rangle, z^1 \langle \mathcal{L} \rangle, \dots, z^{\lfloor \tau n \rfloor - 1} \langle \mathcal{L} \rangle) \leftarrow \text{sel}()$ $\theta_t \leftarrow c[z^k \langle \mathcal{L} \rangle]$ such that $\forall i \in \mathcal{N}_\tau \langle c[z^i \langle \mathcal{L} \rangle] \leq c[z^k \langle \mathcal{L} \rangle] \rangle$
Search for a pds $f$ Estimate each pdf in $\hat{P}_f(\mathcal{Z})$ Create an empty vector of new samples Sample $m$ new samples from $\hat{P}_f(\mathcal{Z})$	$f \leftarrow \text{sea}()$ $\{\hat{P}(\cdot)   \hat{P}(\cdot) \leftarrow \hat{P}_f(\mathcal{Z})\} \leftarrow \text{est}()$ $\mathcal{O} \leftarrow ()$ <b>for</b> $i \leftarrow 0$ <b>to</b> $m - 1$ <b>do</b> $\mathcal{O} \leftarrow \mathcal{O} \sqcup \text{sam}()$
Replace a part of $\mathcal{P}$ with a part of $\mathcal{O}$ Evaluate the new samples in $\mathcal{P}$	$\text{rep}()$ <b>for each</b> unevaluated $\mathcal{P}_i$ <b>do</b> $c[\mathcal{P}_i] \leftarrow C(\mathcal{P}_i)$
Update the generation counter Denote the required iterations by $t_{\text{end}}$	$t \leftarrow t + 1$ $t_{\text{end}} \leftarrow t$

In the IDEA framework, we have that  $\mathcal{N}_\tau = (0, 1, \dots, \tau n - 1)$ ,  $\tau \in [1/n, 1]$ ,  $\text{sel}()$  is the selection operator,  $\text{rep}()$  replaces a subset of  $\mathcal{P}$  with a subset of  $\mathcal{O}$ ,  $\text{ter}()$  is the termination condition,  $\text{sea}()$  is a pds search algorithm,  $\text{est}()$  estimates each pdf and  $\text{sam}()$  generates a single sample from  $\hat{P}_f(\mathcal{Z})$ . The notation  $\hat{P}(\cdot) \leftarrow \hat{P}_f$  means that  $\hat{P}(\cdot)$  is one of the pdfs that is implied by the model  $f$ . The IDEA is a true evolutionary algorithm in the sense that a population of individuals is used from which individuals are selected to generate new offspring with. Using these offspring along with the parent individuals and the current population, a new population is constructed. By referring to the iterations in the IDEA as generations, the evolutionary correspondence is even more obvious. Note that in the IDEA, we have used the approximation notation  $\hat{P}^{\theta_t} f(\mathcal{Z})$  instead of the true distribution  $P^{\theta_t} f(\mathcal{Z})$ . An approximation is required because the determined distribution is based upon samples and the underlying density model is an assumption on the true model. This means that even though we might achieve  $\hat{P}^{\theta_t} f(\mathcal{Z}) = P^{\theta_t} f(\mathcal{Z})$ , in general this is not the case. If we set  $m$  to  $(n - \tau n)$ ,  $\text{sel}()$  to selection by taking the best  $\tau n$  vectors and  $\text{rep}()$  to replacing the worst  $(n - \tau n)$  vectors by the new sampled vectors, we have that  $\theta_{k+1} = \theta_k - \varepsilon$  with  $\varepsilon \geq 0$ . This assures that the search for  $\theta^*$  is conveyed through a monotonically decreasing series  $\theta_0 \geq \theta_1 \geq \dots \geq \theta_{t_{\text{end}}}$ . We call an IDEA with  $m$ ,  $\text{sel}()$  and  $\text{rep}()$  so chosen, a monotonic IDEA. If we set  $m$  in the IDEA to  $n$  and set  $\text{rep}()$  to replace  $\mathcal{P}$  with  $\mathcal{O}$ , we obtain the EDA by Mühlenbein, Mahnig and

Rodriguez [17]. In the EDA however, the threshold  $\theta_t$  cannot be enforced. Note how EDA is thus an instance of IDEA.

## Probability density structure search algorithms

In order to search for a pds, a metric is required that guides the search. In effect, this poses another optimization problem. The metric we use in this paper is a distance metric to the full joint pds  $(\pi^+, \omega^+)$ ,  $\forall i \in L$   $\omega^+ i = i \wedge \pi^+(i) = (i + 1, i + 2, \dots, l - 1)$ . The distance metric is defined by the Kullback–Leibler (KL) divergence. We write  $Y$  instead of  $Z$  from now on to indicate the use of continuous random variables instead of either the discrete or continuous case. Using our definitions, the KL divergence can be written as [7]:

$$D(\hat{P}_{(\pi^+, \omega^+)}(\mathcal{Y}) || \hat{P}_{(\pi, \omega)}(\mathcal{Y})) = -h(\hat{P}_{(\pi^+, \omega^+)}(\mathcal{Y})) + \sum_{i=0}^{l-1} h(\hat{P}(Y_{\omega_i} | Y_{\langle \pi(\omega_i) \rangle})) \quad (3)$$

Let  $a \in L$ ,  $b \in L$  where  $a \in L$  means that  $a$  contains only elements of  $L$ . In equation 3,  $h(Y_{\langle a \rangle})$  is the multivariate differential entropy and  $h(Y_{\langle a \rangle} | Y_{\langle b \rangle})$  is the conditional differential entropy. Let

$$dy_{\langle a \rangle} = \prod_{i=0}^{|a|-1} dy_i$$

be shorthand notation for the multivariate derivative. We then have:

$$h(P(Y_{\langle a \rangle})) = - \int P(Y_{\langle a \rangle})(y_{\langle a \rangle}) \ln(P(Y_{\langle a \rangle})(y_{\langle a \rangle})) dy_{\langle a \rangle} \quad (4)$$

$$h(P(Y_{\langle a \rangle} | Y_{\langle b \rangle})) = h(P(Y_{\langle a \sqcup b \rangle})) - h(P(Y_{\langle b \rangle})) \quad (5)$$

As the term

$$h(\hat{P}_{(\pi^+, \omega^+)}(\mathcal{Y}))$$

in equation 3 is constant, an algorithm that searches for a pds can use the KL divergence by minimizing the sum of the conditional entropies imposed by  $(\pi, \omega)$ . This will cause the pds search algorithm to search for a pds as close as possible to  $(\pi^+, \omega^+)$  subject to additional constraints. The probabilistic models used in previously proposed algorithms range from lower order structures to structures of unbounded complexity. It has been empirically shown by Bosman and Thierens [6] that a higher order pds is required to solve higher order building block problems. We



shortly state three previously introduced pds search algorithms that we use in our experiments. In the univariate distribution, all variables are regarded independently of each other. The PBIL by Baluja and Caruana [2], the cGA by Harik, Lobo and Goldberg [13], the UMDA by Mühlenbein and Paaß [18], and all known approaches in the continuous case prior to the IDEA [10, 22, 23], use this pds. It can be modelled by  $\forall i \in L \quad \pi(i) = () \wedge \omega_i = i$ , giving:

$$\hat{P}_{(\pi, \omega)}(\mathcal{Z}) = \prod_{i=0}^{l-1} \hat{P}(Z_i).$$

In the MIMIC algorithm by De Bonet, Isbell and Viola [4], the pds is a chain which is constrained to  $\pi(\omega_{l-1}) = () \wedge \forall i \in L - (l-1) \quad \pi(\omega_i) = (\omega_{i+1})$ , giving

$$\hat{P}_{(\pi, \omega)}(\mathcal{Z}) = (\prod_{i=0}^{l-2} \hat{P}(Z_{\omega_i} | Z_{\omega_{i+1}})) \hat{P}(Z_{\omega_{l-1}}).$$

To find the chain, an  $O(l^2)$  greedy approximation algorithm is used in MIMIC to minimize the KL divergence. If the pds is constrained so that in addition to having an acyclic pds graph, each node may have at most  $\kappa$  parents, the pds is constrained to  $\forall i \in L \quad |\pi(i)| \leq \kappa$ . This general approach is used in the BOA by Pelikan, Goldberg and Cantú-Paz [19], as well as the LFDA by Mühlenbein and Mahnig [16] and the EBNA by Larrañaga, Etxeberria, Lozano and Peña. In the case of  $\kappa = 1$ , a polynomial time algorithm can be used to minimize the KL divergence [5]. In the case of  $\kappa > 1$ , a greedy algorithm is used that iteratively adds arcs to the pds graph. There are other special case algorithms, such as the optimal dependency trees approach by Baluja and Davies [3] and the ECGA by Harik [12]. Like the LFDA, the ECGA uses minimum description length as a search metric. This metric has the advantage that the resulting pds will not be overly complex. Using the KL divergence, this can only be influenced by adjusting  $\kappa$  because the KL divergence is merely a distance measure from a certain pds to  $(\pi^+, \omega^+)$ . We only regard the three described search algorithms in combination with the KL divergence metric. Note that using the KL metric and the  $(\pi^+, \omega^+)$  pds is merely an instance of the IDEA framework. This is also the case for using a certain pdf. The framework is to be seen separately from the algorithms that can be modelled by it.

## Probability density functions

Next to the pds search algorithms from section 3, we require to specify a pdf to use. It follows from sections 2 and 3 that we require to know the multivariate differential entropy as well as the conditional pdf. In this section, we specify two well known pdfs that we use in our experiments within the IDEA framework. A widely used parametric pdf is the normal pdf. Let

$$\mathcal{S} = (y^0, y^1, \dots, y^{|\mathcal{S}|-1})$$



be the set of selected samples. The sample average in dimension  $j$  is then

$$\bar{Y}_j = \frac{1}{|\mathcal{S}|} \sum_{i=0}^{|\mathcal{S}|-1} \mathbf{y}^i_j.$$

The sample covariance matrix over variables  $\mathbf{y} \langle \mathbf{a} \rangle$  is

$$\frac{1}{|\mathcal{S}|} \sum_{i=0}^{|\mathcal{S}|-1} (\mathbf{y}^i \langle \mathbf{a} \rangle - \bar{\mathbf{Y}} \langle \mathbf{a} \rangle)(\mathbf{y}^i \langle \mathbf{a} \rangle - \bar{\mathbf{Y}} \langle \mathbf{a} \rangle)^T.$$

Let

$$s'_{ij} = \mathbf{S}^{-1}(i, j).$$

The conditional pdf and the entropy can be stated as follows [5]:

$$f_{\mathcal{N}}(y_{\mathbf{a}_0} | \mathbf{y} \langle \mathbf{a} - \mathbf{a}_0 \rangle) = \frac{1}{\sigma \sqrt{2\pi}} e^{\frac{-(y_{\mathbf{a}_0} - \mu)^2}{2\sigma^2}} \quad (6)$$

$$\text{where } \sigma = \frac{1}{\sqrt{s'_{00}}}, \quad \mu = \frac{\bar{Y}_{\mathbf{a}_0} s'_{00} - \sum_{i=1}^{|\mathbf{a}|-1} (y_{\mathbf{a}_i} - \bar{Y}_{\mathbf{a}_i}) s'_{i0}}{s'_{00}}$$

$$h(Y \langle \mathbf{a} \rangle) = \frac{1}{2} (|\mathbf{a}| + \ln((2\pi)^{|\mathbf{a}|} \det(\mathbf{S}))) \quad (7)$$

The non-parametric normal kernels pdf places a normal pdf over every available sample point. Let  $s_i$  be a fixed standard deviation in the  $i$ -th dimension. The conditional pdf and the entropy can then be stated as follows [8]:

$$f_{\mathcal{N}_K}(y_{\mathbf{a}_0} | \mathbf{y} \langle \mathbf{a} - \mathbf{a}_0 \rangle) = \sum_{i=0}^{|\mathcal{S}|-1} \nu_i \frac{1}{s_{\mathbf{a}_0} \sqrt{2\pi}} e^{\frac{-(y_{\mathbf{a}_0} - \mathbf{y}^i_{\mathbf{a}_0})^2}{2s_{\mathbf{a}_0}^2}} \quad (8)$$

$$\text{where } \nu_i = \frac{e^{-\sum_{j=1}^{|\mathbf{a}|-1} \frac{(y_{\mathbf{a}_j} - \mathbf{y}^i_{\mathbf{a}_j})^2}{2s_{\mathbf{a}_j}^2}}}{\sum_{k=0}^{|\mathcal{S}|-1} e^{-\sum_{j=1}^{|\mathbf{a}|-1} \frac{(y_{\mathbf{a}_j} - \mathbf{y}^k_{\mathbf{a}_j})^2}{2s_{\mathbf{a}_j}^2}}}$$

$$h(Y \langle \mathbf{a} \rangle) = \quad (9)$$

$$\frac{1}{2} \ln \left( |\mathcal{S}|^2 (2\pi)^{|\mathbf{a}|} \prod_{j=0}^{|\mathbf{a}|-1} s_{\mathbf{a}_j}^2 \right) - \int f_{\mathcal{N}_K}(\mathbf{y} \langle \mathbf{a} \rangle) \ln \left( \sum_{i=0}^{|\mathcal{S}|-1} e^{-\sum_{j=0}^{|\mathbf{a}|-1} \frac{(y_{\mathbf{a}_j} - \mathbf{y}^i_{\mathbf{a}_j})^2}{2s_{\mathbf{a}_j}^2}} \right) d\mathbf{y} \langle \mathbf{a} \rangle$$

An alternative pdf to the two described above, is the histogram pdf. Using this pdf however does not scale up very well [7] and leads to an exponential iteration running time in the order of  $rk$  where  $r$  is the amount of bins to use in each dimension. The normal pdf is very efficient but very cluster insensitive. The normal kernels pdf is very sensitive to clusters but may very quickly overfit the data. In addition, the running time each iteration for using the latter pdf tends to be a lot greater than for the normal pdf.

## Conclusion

In this project, we have applied the search for good probability density models to continuous spaces. We have formalised the notion of building and using probabilistic models in a new framework named IDEA. We have shown how we can adjust existing techniques to be used in the continuous case.

We have utilised the algorithmic system IDEA for demonstrating iterated density estimation evolutionary algorithms. These calculations utilise density estimation procedures to probability distribution over the factors that code an issue to perform optimisation. To this end, a probability density structure should be found and therefore be utilised in density estimation. For a bunch of existing search algorithms, we have applied and tried them in the IDEA framework using two distinctive density estimation models.. The experiment demonstrates that building and utilising probabilistic models for continuous optimisation problems is promising. This, in blend with its demonstrating abilities, shows that the IDEA is general, appropriate and compelling.

## References

1. [https://link.springer.com/chapter/10.1007/3-540-45356-3\\_75](https://link.springer.com/chapter/10.1007/3-540-45356-3_75)
2. [https://homepages.cwi.nl/~bosman/publications/2000\\_expandingfromdiscr etc.pdf](https://homepages.cwi.nl/~bosman/publications/2000_expandingfromdiscr etc.pdf)