

# Large scale estimation of distribution algorithms with adaptive heavy tailed random projection ensembles

Sanyang, Momodou; Kaban, Ata

*License:*

Other (please specify with Rights Statement)

*Document Version*

Peer reviewed version

*Citation for published version (Harvard):*

Sanyang, M & Kaban, A 2019, 'Large scale estimation of distribution algorithms with adaptive heavy tailed random projection ensembles', *Journal of Computer Science and Technology*, vol. 34, no. 6, pp. 1241-1257. <<http://jcst.ict.ac.cn/EN/10.1007/s11390-019-1973-1>>

[Link to publication on Research at Birmingham portal](#)

**Publisher Rights Statement:**

This is a post-peer-review, pre-copyedit version of an article published in *Journal of Computer Science and Technology*. The final authenticated version is available online at: <http://dx.doi.org/10.1007/s11390-019-1973-1>

**General rights**

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

**Take down policy**

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

# Large Scale Estimation of Distribution Algorithms with Adaptive Heavy Tailed Random Projection Ensembles

Momodou L. Sanyang<sup>1,2</sup>, Ata Kabán<sup>1</sup>, Member IEEE

<sup>1</sup>*School of Computer Science, University of Birmingham, Edgbaston, B15 2TT, Birmingham, UK.*

<sup>2</sup>*School of Information Technology and Communication, University of The Gambia, MDI Road, Kanifing, P.O. Box 3530, Serekunda, The Gambia.*

E-mail: MLSanyang123@gmail.com, MLSanyang@utg.edu.gm ; A.Kaban@cs.bham.ac.uk

Received xxxxx xx 2018; revised xxxxx xx, 2019

**Abstract** We present new variants of Estimation of Distribution Algorithms (EDA) for large scale continuous optimisation, that extend and enhance a recently proposed random projection (RP) ensemble based approach. The main novelty here is to depart from the theory of RPs that required (sub-)Gaussian random matrices for norm-preservation, and instead for the purposes of high dimensional search we propose to employ random matrices with independent and identically distributed entries drawn from a t-distribution. We analytically show that the implicitly resulting high dimensional covariance of the search distribution is enlarged as a result. Moreover, the extent of this enlargement is controlled by a single parameter, the degree of freedom. For this reason, in the context of optimisation, such heavy tailed random matrices turn out to be preferable over the previously employed (sub-)Gaussians. Based on this observation, we then propose novel covariance adaptation schemes that are able to adapt the degree of freedom parameter during the search, and give rise to a flexible approach to balance exploration versus exploitation. We perform a thorough experimental study on high dimensional benchmark functions, and provide statistical analyses that demonstrate state-of-the-art performance of our approach when compared with previously existing alternatives in problems with 1000 search variables.

**Keywords** covariance adaptation, estimation of distribution algorithm, random projection ensemble, t-distribution

## 1 Introduction

Optimisation over high dimensional search spaces is a key task in many modern applications, including scientific, engineering and management problems. Estimation of Distribution Algorithms (EDAs) are black-box optimisation methods based on probabilistic modelling, which are becoming increasingly important in the Evolutionary Computation (EC) community [1]. EDAs avoid the usage of arbitrary operators of traditional EC approaches, instead they estimate a probabilistic model from the fittest individuals, and advance the search for the global optimum by sampling the esti-

mated model. In this way, EDAs will not only simplify the design of an optimisation method, but also allow the algorithms to learn the structure of the search space and guide the search in further promising directions [2].

However, the performance of most existing EDAs deteriorates significantly in high dimensional search spaces. This is because EDA requires model building from empirical data, which is susceptible to the curse of dimensionality [3]. Unless an enormous budget of fitness evaluations is available, the sample size is insufficient for reliable model estimation. Consequently the structure of the problem is often severely mis-estimated

from limited samples – see [4, 5] for examples. In turn, a badly estimated model typically results in premature convergence [5].

To mitigate the curse of dimensionality, simplified models such as univariate models, or limited dependency models are often preferred instead of a full multivariate model in practice. The simplest of such approach is described in [6], called UMDAc. The model building in this univariate approach is simply done under the independence assumption [6]:  $P(\mathbf{x}) = \prod_{i=1}^d P(x_i)$ , where  $d$  is the dimensionality of the search space. More refined approaches such as the sep-CMA-ES [7] and the univariate AMaLGaM [8] use only diagonal elements of a rotated full covariance matrix. However, the univariate models assume that all the variables are independent in some coordinate basis, hence they may be expected to perform well on problems that have no dependency between variables (i.e., separable) in that basis.

Non-separable problems are those in which the above assumption does not hold – that is, the design variables have multiple inter-dependencies. The approach known as MIMIC, proposed by De Bonet et al [9] sets out to capture bivariate interactions between decision variables by sampling from the pairwise joint distribution between variables. Despite that MIMIC is able to outperform univariate models, the majority of optimisation problems will have larger groups of interacting design variables. Several proposals have been put forth to explicitly capture multivariate dependencies by building graphical dependency networks, for example Bayesian Networks [10]. But, from statistical, computational and memory points of view, learning probabilistic graphical models from the sample of fit individuals is highly expensive and requires a large sample [11]. Thus, the scaling-up of these model building processes to high dimensional problems remains chal-

lenging.

### 1.1 Related Work In Large-scale Black-box Optimisation

There have been many efforts to alleviate the problems of high dimensional search, resulting in several algorithms being proposed recently, including EDA-type methods. Here we limit ourselves to a few that are most relevant to the present work. Approaches can be grouped into cooperative co-evolution based methods, evolutionary strategies, latent variable models, hybrid methods and compression based divide and conquer.

One of the first cooperative co-evolution (CC) methods is CC with Variable Interaction Learning (CCVIL) proposed by Weicker et al. in [12]. It is a deterministic method to uncover dependencies between decision variables. A more recent model-based formulation of a similar idea is EDA with Model Complexity Control (EDA-MCC) [13], which also employs a deterministic algorithm; it groups variables into two disjoint subsets – those that are weakly correlated with other variables, and those that are strongly correlated. The strength of correlations is estimated relative to a threshold.

Multilevel Cooperative Co-evolution (MLCC) proposed in [14] groups the decision variables of a problem randomly to tackle problems in high-dimensional optimisation. The groups are then optimised jointly, but separately from other groups.

MA-SW-Chain [15] is a hybrid algorithm that assigns local search intensity to each individual by chaining dissimilar local search applications. It is an extension of the MA-CMA-Chain algorithm for high dimensional regime, and was the winner of the CEC2010-2012 large scale global optimisation competition on 1000 dimensional benchmark problems.

Another related method employs latent variable models to sample new individuals with low-dimensional

latent vectors. Shin et al [16] proposed a latent variable model such as the Helmholtz machine and probabilistic principal component analysis to estimate the probability distribution of given data. The model considers latent variables which have much lower dimension than input variables for optimising a high dimensional function. In doing so, the method will mitigate the curse of dimensionality. A different way to exploit hidden low dimensional structure was proposed in [17], suggesting that a notion of intrinsic dimension of a function can replace the ambient dimension of the inputs to the function, provided that the function has a special structure. However, in the absence of this structure this method is unsuitable.

A recent approach is the random matrix theory based EDA in [4], which devises a compression-based divide and conquer strategy. The idea is to project the high dimensional selected individuals to many independent low dimensional random subspaces, carry out the costly sampling and estimation operations concurrently in these subspaces, and combine the new samples to form the next generation. Such ensemble can exhibit a strong smoothing effect that facilitates estimation from small samples.

## 1.2 Contributions

In this paper, we further enhance the random projection (RP) ensemble approach of [4], by employing a parametrised family of heavy-tailed random matrices to replace the more conventional sub-Gaussian ones. Although this results in a minor change in the implementation, as we shall see, it allows for a better exploration of high dimensional search spaces. An early version with promising empirical results of this idea appeared in [18]\*.

The increased exploration abilities of heavy tailed

distributions are well known in evolutionary search, starting from pioneering work by Yao et al. [19] in the univariate setting, and more recent work in the multivariate EDA framework – see e.g. [20–22] and references therein. However, in the regime of large search spaces of dimensionality well beyond 100 variables, the use of such heavy tailed distributions in EDA is not straightforward – as demonstrated in [22], a heavy tailed search distribution becomes increasingly counter-productive as it loses sight of the direction of the search. Here we avoid these problems as we employ heavy tailed distributions in a combination of RPs rather than directly in the role of a search distribution.

Our approach is based on a detailed analysis of the effect of the distribution of entries of random matrices used in the algorithm in terms of the resulting *aggregated* covariance of the multivariate search distribution (which is not heavy tailed but Gaussian, as we shall see later). We should also note that this type of analysis is entirely different, and complementary with analyses of the dynamics or the running time, as the latter is currently feasible only in univariate models on some very specific problems [23]. The purpose of our analysis in this work is (i) to shed light on some limitations of borrowing existing tools from the area of RPs as done in previous work, and (ii) to construct novel algorithm variants that bypass these limitations.

In particular, our analysis reveals that, while the use of RPs is useful for dealing with high dimensional optimisation problems, the existing RP theory is not well aligned with the needs of optimisation. More precisely, the use of sub-Gaussian RP matrices has been originally borrowed from a literature that aims at approximately preserving Euclidean geometry. This goal is different from ours and, as we shall see, adopting the same approach falls short of exploration ability for

---

\*Received a Runner-Up Student Paper award.

search. In turn, the use of heavy tailed random matrices is very unusual in the context of the RP literature, as they do not have good distance preservation properties, but they will turn out to have advantages for high dimensional search.

The specific contributions of this paper are as follows:

- We analytically show that the excess kurtosis of the entries of the RP matrices in the ensemble has a precise role in the aggregated covariance of the ensemble. The latter is the covariance of the high dimensional search distribution, therefore this analysis reveals the following: A distribution with negative excess kurtosis will shrink the aggregated covariance in comparison with that obtained when using the Gaussian distribution, and a distribution with positive excess kurtosis will enlarge it.
- We then propose to employ random matrices with i.i.d. t-distributed entries – this subsumes the Gaussian ensemble of [4] as a special case, when the degree of freedom is infinite. However, the degree of freedom parameter offers a means to enlarge the aggregated covariance in favour of better exploration in high dimensions. In particular, we show that the extent of this enlargement is a decreasing function of the degree of freedom parameter.
- Based on this, we then develop simple methods to adapt the degree of freedom parameter during the search, which essentially implement novel covariance adaptation schemes specifically for high dimensional search.
- We present a thorough experimental study on 1000-dimensional multi-modal test functions from

the CEC2010-2012 large scale global optimisation competition benchmark, and demonstrate superior or state of the art performance in statistical comparisons with a number of existing alternatives.

## 2 Heavy-tailed Random Matrices For Continuous EDA Optimisation

Random Projections (RP) have already been used with success in large scale continuous EDA [4]. However the theory of RP in the literature is aimed at preservation of the Euclidean distances and norms. The goal in optimisation is different: it is more important to have a good exploration-exploitation tradeoff.

With this goal in mind, in this section we shall present a RP ensemble based large scale multivariate Gaussian EDA that employs random matrices with heavy tailed entries drawn i.i.d. from the family of t-distributions. We shall present this algorithm first, to fix ideas, so that we see how these random matrices enter into the search algorithm. We will then analytically justify our choice for the family of t-distributions afterwards in the subsequent subsection.

### 2.1 The tRP-Ens-EDA algorithm: A basic variant

Our algorithm is a modification of the random projection ensemble based EDA (RP-Ens-EDA) of [4], and we refer to our new variant as tRP-Ens-EDA. The pseudo-code of tRP-Ens-EDA is given in Algorithm 1.

The tRP-Ens-EDA proceeds by initially generating a population of individuals randomly everywhere in the search space, and selects the  $T$  fittest ones based on their fitness values. This is the set  $\mathcal{P}^{\text{fit}}$  in Algorithm 1.  $M$  independent random projection (RP) matrices  $\mathbf{R}_i$  are then generated, each defines a random subspace of dimension  $k \ll d$ , where  $d$  is the dimension of the search space. These

are created in order to project the fittest individuals onto low dimensional subspaces defined by them.

---

**Algorithm 1:** Algorithm with entries of  $\mathbf{R}$  from  $t$ -distribution (tRP-Ens-EDA)

---

**Input:**  $k, M, N, MaxFE$

**Output:** Best individual  $\mathbf{x}^*$ ; Best fitness  $f^*$

```

1  $\mathcal{P} \leftarrow$  Initialise a population uniformly at random
  in the search space.
2  $\nu \leftarrow$  Initialise a degree of freedom
3 while  $MaxFE > 0$  do
4    $MaxFE \leftarrow MaxFE - N$ .
5    $\mathbf{f} \leftarrow$  Evaluate the fitnesses of  $\mathcal{P}$ .
6    $\mathcal{P}^{fit} \leftarrow$  Select the  $T$  individuals of  $\mathcal{P}$  with
  best fitness using truncation selection.
7    $\mathbf{x}^* \leftarrow$  the individual with the optimal fitness
  so far.
8    $f^* \leftarrow$  the optimal fitness value so far.
9    $\nu \leftarrow$  Decide a new  $\nu$  with an adaptive
  method.
10   $\boldsymbol{\mu} \leftarrow$  Estimate the mean of  $\mathcal{P}^{fit}$  using
  maximum likelihood estimation (MLE).
11   $\mathcal{P}^{fit} \leftarrow \mathcal{P}^{fit} - \boldsymbol{\mu}$ .
12   $\{\mathbf{R}_i\}_{i=1:M} \leftarrow$  Generate  $M$  independent
  random matrices with entries drawn i.i.d.
  from a  $t$ -distribution with mean 0 and
  variance  $\frac{1}{d}$ .
13  foreach  $i \in \{1, 2, \dots, M\}$  do
14     $\mathcal{Y}^{\mathbf{R}_i} \leftarrow \mathbf{R}_i \mathcal{P}^{fit}$  projects  $\mathcal{P}^{fit}$  to a random
     $k$ -dimensional subspace.
15     $\boldsymbol{\Sigma}^{\mathbf{R}_i} \leftarrow$  Estimate the  $k \times k$  covariance
    matrix of  $\mathcal{Y}^{\mathbf{R}_i}$  using MLE.
16    Sample  $N$  new points  $\mathbf{y}_1^{\mathbf{R}_i}, \dots, \mathbf{y}_N^{\mathbf{R}_i}$  i.i.d.
    from  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{\mathbf{R}_i})$ .
17   $\mathcal{P}^{new} \leftarrow$ 
     $\sqrt{\frac{dM}{k}} [\frac{1}{M} \boldsymbol{\Sigma}_{i=1}^M \mathbf{R}_i^T \mathbf{y}_1^{\mathbf{R}_i}, \dots, \frac{1}{M} \boldsymbol{\Sigma}_{i=1}^M \mathbf{R}_i^T \mathbf{y}_N^{\mathbf{R}_i}] + \boldsymbol{\mu}$ .
18   $\mathcal{P} \leftarrow \mathcal{P}^{new}$ .
19  Replace an individual in  $\mathcal{P}$  with  $\mathbf{x}^*$  (elitism).
```

---

The number of such subspaces,  $M$ , and their dimension,  $k$ , are parameters of the method along with the population size  $N$ . For all of these parameters we will use the default values determined in [4]. i.e.,  $M = \lceil \frac{3d}{k} \rceil$ ,  $k = 3$ ,  $N = 300$ , and  $|\mathcal{P}^{fit}| = N/4$ , as the context and meaning of these parameters is the same as in that work. Another input parameter is the maximum fitness evaluations allowed, MaxFE, which will be varied depending on the goal of the experiments.

Once the set  $\mathcal{P}^{fit}$  is determined, its sample mean is estimated in step 10 and used to center the points in step 11. The entries of the RP matrices  $\mathbf{R}_i$  are drawn i.i.d. from a  $t$ -distribution with mean 0 and variance  $\frac{1}{d}$ . This is the only difference from the original RP-Ens-EDA [4]. It is implemented by sampling from  $t(0, 1, \nu)$  first, and then multiplying the samples by  $\sqrt{\frac{\nu-2}{\nu d}}$  so that their variance becomes  $1/d$ . This choice for the variance is to ensure that we recover the original scale in Step 17 without having to modify the scaling factor of the original RP-Ens-EDA. When  $d$  is large,  $\mathbf{R}_i$  have nearly orthonormal rows, as a result of the concentration of norms in high dimensions, since the coordinates in each row are independent – this is a very generic property of high dimensional probability spaces that both [4] and our method exploit. So, pre-multiplying with  $\mathbf{R}_i$  is almost like orthogonally projecting the points from the  $d$  dimensional space to a  $k$  dimensional subspace, which shortens the lengths of vectors by a factor of  $\sqrt{\frac{k}{d}}$  and the standard deviation gets reduced by a factor of  $\sqrt{M}$  after averaging. Therefore, the scaling factor needed to recover the original scale is  $\sqrt{\frac{dM}{k}}$ .

Step 14 projects the good samples,  $\mathcal{Y}^{\mathbf{R}_i}$  down to the subspaces of dimension  $k$ , then estimates the  $k \times k$  covariance matrices in each subspace and samples  $N$  new points in each subspace using  $k$ -dimensional multivariate Gaussian distributions. Step 17 averages the obtained points from these subspaces to produce the new population  $\mathcal{P}$ . In addition, we will use elitism in practice, so the best fitness individual always survives to the next generation.

## 2.2 Why $t$ -distributed random projections?

Let us restrict our attention to the procedure by which the new generation is created, as this is what distinguishes our algorithm from the vanilla multivariate EDA. We want to analyse the aggregated covariance

that the ensemble of RP creates implicitly, which is in fact the covariance of the high dimensional search distribution that the new generation follows.

Our analysis will allow quite general RP matrices, and only require their entries to be drawn i.i.d. from a 0-mean symmetric distribution having finite first four moments. This includes the Gaussian, all sub-Gaussians, and the family of  $t$ -distributions with degree of freedom at least 5. We take advantage of this generality to develop our arguments and explain our choice for the  $t$ -distributions, which, for readers familiar with the random projection literature may seem very unusual.

We start by computing the covariance of the new population in step 17 of Algorithm 1, conditional on fixing the random projection matrices  $\mathbf{R}_i, i = 1 : M$ . We will then condition on the sample of fit individuals and look at the effect of  $\mathbf{R}_i, i = 1 : M$  by computing the expectation of this ensemble covariance with respect to  $\mathbf{R}_i, i = 1 : M$ . If  $M$  is large enough, this expectation approximates well the finite  $M$  case. Moreover, although not pursued here, tools from [24] may be used to determine the sufficient ensemble size.

**Proposition 1** *Conditionally on  $\mathbf{R}_i, i = 1 : M$ , the new generation produced at Step 17 of Algorithm 1 is i.i.d. Gaussian with mean  $\boldsymbol{\mu}$  and  $d \times d$  covariance of the following form:*

$$\boldsymbol{\Sigma}_{rp} = \frac{d}{kM} \sum_{i=1}^M \mathbf{R}_i^T \mathbf{R}_i \boldsymbol{\Sigma} \mathbf{R}_i^T \mathbf{R}_i, \quad (1)$$

where  $\boldsymbol{\Sigma}$  is the maximum likelihood sample covariance of the original selected individuals in  $\mathcal{P}^{fit}$ .

The proof is deferred to Appendix A.

Next, we want to see the effect of the random  $\mathbf{R}_i$ 's on  $\boldsymbol{\Sigma}_{rp}$ . To this end, we condition on  $\boldsymbol{\Sigma}$ , and look at the expectation  $E_{\mathbf{R}}[\boldsymbol{\Sigma}_{rp}]$ .

We require the following definition.

**Definition 1** *The excess kurtosis of a random variable  $x$  is defined as:*

$$K = \frac{E[x^4]}{E[x^2]^2} - 3.$$

We have the following result.

**Theorem. 1** *Let  $\mathbf{R}$  be a  $k \times d$  random matrix,  $k < d$ , with entries drawn i.i.d. from a symmetric distribution with 0-mean, finite first four moments, and excess kurtosis  $K$ . Let  $\boldsymbol{\Sigma}$  be a  $d \times d$  fixed positive semi-definite matrix with eigenvalues  $\lambda_1, \dots, \lambda_d$ .*

$$E_{\mathbf{R}}[\boldsymbol{\Sigma}_{rp}] = \frac{1}{d} \left[ (k+1)\boldsymbol{\Sigma} + \text{Tr}(\boldsymbol{\Sigma})\mathbf{I}_d + K \sum_{i=1}^d \lambda_i \mathbf{A}_i \right], \quad (2)$$

where  $\mathbf{A}_i$  are  $d \times d$  diagonal matrices with their  $j$ -th diagonal elements being  $\sum_{a=1}^d U_{ai}^2 U_{aj}^2$  and  $U_{ai}$  is the  $a$ -th entry of the  $i$ -th eigenvector of  $\boldsymbol{\Sigma}$ .

*Proof* From eq. (1), the expectation of  $\boldsymbol{\Sigma}_{rp}$  is  $\frac{d}{k} E_{\mathbf{R}}[\mathbf{R}^T \mathbf{R} \boldsymbol{\Sigma} \mathbf{R}^T \mathbf{R}]$ , which we compute using a result from [25], originally derived in a very different context (namely to analyse compressive least square regression).

Under the condition on  $\mathbf{R}$  stated in Theorem 1, Lemma 2 from [25] proved that  $E[\mathbf{R}^T \mathbf{R} \boldsymbol{\Sigma} \mathbf{R}^T \mathbf{R}] =$

$$k \cdot E[R_{i,j}^2]^2 [(k+1)\boldsymbol{\Sigma} + \text{Tr}(\boldsymbol{\Sigma})\mathbf{I}_d + K \sum_{i=1}^d \lambda_i \mathbf{A}_i], \quad (3)$$

where  $\mathbf{A}_i$  are the diagonal matrices defined as in Theorem 1, and  $R_{i,j}$  is a generic entry of matrix  $\mathbf{R}$ .

Now, by design, we have  $E[R_{i,j}^2] = \frac{1}{d}$ . Replacing eq. (3) into eq. (1) completes the proof.  $\square$

Observe that, from the definition of the matrices  $\mathbf{A}_i$ , and since  $\boldsymbol{\Sigma}$  is positive semi-definite,  $\lambda_i$  are non-negative – consequently the sum in the last term in the brackets in eq. (2), that is  $\sum_{i=1}^d \lambda_i \mathbf{A}_i$  is always a diagonal matrix with non-negative diagonal entries.

Now, let us look at the multiplier of this term,  $K$ . In previous work [4] the entries of  $\mathbf{R}$  were Gaussian or sub-Gaussian. In the Gaussian case,  $K = 0$ . Hence in that case the last term cancels out:  $K \sum_{i=1}^d \lambda_i \mathbf{A}_i = 0$ . Nevertheless, what remains is a regularised version of the sample covariance estimate  $\mathbf{\Sigma}$ . This is always non-singular, so the search does not get stuck in a subspace of the search space.

However, if the entries of  $\mathbf{R}$  are drawn from a sub-Gaussian distribution, we find that certain choices may be a bad idea – it depends on the excess kurtosis of the distribution.

Indeed, take the Rademacher distribution below, and note this is sub-Gaussian:

$$R_{ij} \sim_{i.i.d.} \begin{cases} 1 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/2. \end{cases}$$

Random matrices with i.i.d. entries from this distribution have been originally proposed by [26] as an implementation-friendly alternative to the Gaussian RP, and it was subsequently considered for RP-EDA-Ensemble optimisation in [4].

It is easy to check that for this distribution we have  $E[R_{ij}^4] = 1$ , and  $E[R_{ij}^2] = 1$ . Therefore, the excess kurtosis is  $K = 1 - 3 = -2$ . Consequently,  $K \sum_{i=1}^d \lambda_i \mathbf{A}_i = -2 \sum_{i=1}^d \lambda_i \mathbf{A}_i$ . This is a diagonal matrix with all non-positive diagonal entries. Therefore this term diminishes the regularisation effect of the term  $Tr(\mathbf{\Sigma})\mathbf{I}_d$  in eq. (2). It also shrinks the ensemble covariance matrix in comparison with that of a Gaussian RP ensemble.

We can go even further, and construct a simple example where the term  $K \sum_{i=1}^d \lambda_i \mathbf{A}_i$  cancels out the effect of  $Tr(\mathbf{\Sigma})\mathbf{I}_d$  completely in some directions of the search space, so we are left with maximum likelihood estimates in those directions that may be not sufficiently accurate due to the small number of high-fitness individuals relative to the high dimensionality of the search

space. We give such an example in Appendix B.

Although such examples may seem contrived, the shrinking effect of distributions with negative kurtosis is undesirable for search, and instead we would like to use the insights of this analysis to come up with a better alternative.

If  $K > 0$ , then the term,  $K \sum_{i=1}^d \lambda_i \mathbf{A}_i$  is a diagonal positive semi-definite matrix, hence it always adds a non-negative quantity to the diagonals of the ensemble covariance. Hence, our idea is to use a family of distributions with positive kurtosis, and moreover, adapt  $K$  during the search. Based on the discussion above, this would enlarge the covariance  $E_{\mathbf{R}}[\mathbf{\Sigma}_{rp}]$  adaptively to the necessary extent. Covariance adaptation has been a successful technique in the EDA literature [27], and our idea makes it feasible to implement it in new ways, specifically for large scale EDA search through a more suitable random projection ensemble.

To achieve this, we have chosen the family of  $t$ -distributions with degree of freedom at least 5. Indeed, for this family of distributions the following holds.

**Proposition 2** *The excess kurtosis of  $t(0, 1, \nu)$  distribution with degree of freedom  $\nu$ , is:*

$$K = \frac{6}{\nu - 4}, \text{ provided } \nu > 4.$$

The proof is deferred to the Appendix C.

**Corollary 1 (to Proposition 2)** *Changing the variance leaves the excess kurtosis unchanged.*

The proof is immediate and given in Appendix C for completeness.

We replace this into Theorem 1, and obtain the following.

**Theorem. 2** *If  $\mathbf{R}$  is sampled i.i.d. from a  $t$ -distribution with degree of freedom  $\nu \geq 5$ , mean 0 and*

variance  $1/d$ , then the ensemble covariance for the high dimensional search distribution is:

$$E_{\mathbf{R}}[\boldsymbol{\Sigma}_{rp}] = \frac{1}{d} \left[ (k+1)\boldsymbol{\Sigma} + \text{Tr}(\boldsymbol{\Sigma})\mathbf{I}_d + \frac{6}{\nu-4} \sum_{i=1}^d \lambda_i \mathbf{A}_i \right], \quad (4)$$

where the matrices  $\mathbf{A}_i$  are the same as defined in Theorem 1.

The Gaussian RP corresponds to  $\nu \rightarrow \infty$ , in which case the last term in the bracket in eq. (4) vanishes:  $\lim_{\nu \rightarrow \infty} \frac{6}{\nu-4} \sum_{i=1}^d \lambda_i \mathbf{A}_i = 0$ . In turn, by any finite choice of the degree of freedom parameter for the  $t$ -distribution we will be adding a positive semi-definite matrix to this covariance – this makes it larger and gives it more chance to explore the search space.

### 3 Adaptive Degree Of Freedom Parameter

In this section we devise methods to adapt the parameter  $\nu$  during the search.

Parameter setting methods are dichotomised into tuning and controlling [28]. Tuning means finding a good value by trial and error before running the algorithm and then fixing this value throughout the evolutionary process. While it may seem convenient, the trial-and-error phase consumes part of the budget of function evaluations. On the other hand, parameter control starts with an initial value which is then updated during the search, based on partial outcomes of the algorithm [28]. In other words, the latter type of methods try to adapt the control parameters automatically in order to adjust the algorithm to the problem during the search [29]. Hence these type of methods need us to devise an algorithm that acts as an adaptation strategy. This is what we pursue in the remainder of this section.

In the sequel we present three different adaptation schemes. In each of these, the parameter we try to control is the degree of freedom of the  $t$ -distributed entries

of our random projection matrices ( $\mathbf{R}_i, i = 1 : M$ ).

To simplify the computations we shall approximate  $\boldsymbol{\Sigma}_{rp}$  by an upper bound on it in positive semi-definite ordering, which is independent of the eigenvectors of  $\boldsymbol{\Sigma}$ . Observe that  $\sum_{i=1}^d \lambda_i \mathbf{A}_i \preceq \text{Tr}(\boldsymbol{\Sigma}) \cdot \mathbf{I}_d$ , so we have:

$$E_{\mathbf{R}}[\boldsymbol{\Sigma}_{rp}] \preceq \frac{1}{d} \left[ (k+1)\boldsymbol{\Sigma} + \text{Tr}(\boldsymbol{\Sigma})\mathbf{I}_d \left( 1 + \frac{6}{\nu-4} \right) \mathbf{I}_d \right].$$

#### 3.1 Adaptive Variance Scaling based method (AVS)

An Adaptive Variance Scaling (AVS) method was originally proposed for EDA by Bosman [27]. In each generation, it decides how to update a parameter for the next generation based on the *difference in the best fitness* between the latest two consecutive generations. The original approach [27] is to scale the sample covariance matrix  $\boldsymbol{\Sigma}$  multiplicatively, by a value  $c \in [1, 10]$ , where  $c$  is adapted (i.e., either increased or decreased) by a coefficient  $\eta \in (0, 1)$  as follows. If the best fitness improved then  $c$  is increased as  $c/\eta$  in order to explore, otherwise  $c$  is decreased as  $c\eta$  for exploitation. The authors of [27] put  $c$  to 10 if it becomes greater than 10 or smaller than 1, in order to stimulate exploration while preventing a blow-up of the covariance.

Note, in our context of RP ensemble based algorithm, we do not need to artificially introduce the  $c$  parameter as the original AVS did, as we can use the degree of freedom parameter  $\nu$  to emulate the updating heuristic of the original AVS. The pseudo-code of this

adaptation scheme is given in Algorithm 2.

---

**Algorithm 2:** Adaptive variance based method (AVS)

---

**Input:** Index of current generation  $t$ ; Degree of freedom  $\nu^t$ ; Best fitness  $b^t, b^{t-1}$ ; Coefficient  $\eta$

**Output:** Degree of freedom  $\nu^{t+1}$

```

1  $c \leftarrow 1 + \frac{6}{\nu^{t-4}}$ .
2 if  $b^t$  is no better than  $b^{t-1}$  then
3   |  $c \leftarrow c\eta$ . //  $c$  gets smaller,  $\nu$  gets larger
4 else
5   |  $c \leftarrow c/\eta$ . //  $c$  gets larger,  $\nu$  gets smaller
6  $\nu^{t+1} \leftarrow \frac{6}{c-1} + 4$ .
7 if  $\nu^{t+1} < 5$  or  $\nu^{t+1} > 124$  then
8   |  $\nu^{t+1} \leftarrow 5$ . //to stimulate exploration
```

---

From eq. (4) we have seen that smaller values for  $\nu$  allow for exploration and larger  $\nu$ -values encourage exploitation in tRP-Ens-EDA. Also observe that  $\nu$  has a non-linear effect on the extent of enlarging the covariance of the new population. We set:

$$c = 1 + 6/(\nu - 4) \quad (5)$$

which ensures that  $\nu \geq 5$  as required – so that we can directly apply the same adaptation rules as the original AVS approach of [27], as follows. Suppose the tRP-Ens-EDA is at its  $t$ -th generation, and denote by  $b^t$  and  $b^{t-1}$  the best fitness values of the latest two generations respectively. Recall, tRP-Ens-EDA uses elitism, so we will never observe a worse-than previous value of the best fitness. Therefore, our AVS will decrease  $c$  (hence increase  $\nu$ ) for exploitation if  $b^t$  is no better than  $b^{t-1}$  before the elitism is applied, and increase  $c$  (hence decrease  $\nu$ ) for exploration if  $b^t$  is strictly better than  $b^{t-1}$ . This adaptation strategy is exactly the same as that originally proposed in [27], and the only difference is that it is now built into our RP-ensemble based EDA for large scale searches.

Finally, from eq. (5) we have:

$$\nu = 6/(c - 1) + 4 \quad (6)$$

and, as in the original method of [27], we will also con-

strain the range of  $\nu \in [5, 124]$ , which corresponds to  $c \in [1.05, 7]$ . If  $\nu$  goes beyond this range, then it is put to 5 to trigger exploration, which is the same strategy as in [27].

### 3.2 A 1/5-success rule based method

A classic alternative adaptation scheme is the well-known 1/5-success rule (or 1/5 rule). It is one of the oldest adaptation heuristics, and it is widely used in Evolutionary Strategies (ES) [30]. The 1/5 rule does not look at the improvement in the best fitness, but instead it looks at the overall fitness improvement of the population. It measures this as the following success probability:

$$P_s = \frac{\#\text{individuals with improved fitness}}{\text{Population size}}.$$

If  $P_s > 1/5$  then it is said that an improvement is detected. In our case, when such improvement is detected then we decrease  $\nu$  for exploration, otherwise we increase  $\nu$  in order to exploit – in the same manner as we did in AVS in section 3.1.

However, in ESs it is easy to compute  $P_s$ , because the individuals can be compared with their parents after a mutation – but in EDA there is no explicit connection between a particular individual and its descendant in the next generation, as the new individuals are sampled from a probability distribution that the whole previous generation contributed to. To apply the 1/5 rule in EDAs we use the following heuristic. At the  $t$ -th generation, given the fitness values of the population with  $N$  individuals  $\mathbf{f}^t$ , and the fitness values of the individuals of the previous generation  $\mathbf{f}^{t-1}$ , we sort both  $\mathbf{f}^t$  and  $\mathbf{f}^{t-1}$ , and compare between the corresponding elements  $f_i^t$  and  $f_i^{t-1}$ , where  $i \in \{1, \dots, N\}$ . This way we count the number of fitness values  $f_i^t$  that are better than the corresponding  $f_i^{t-1}$ , and use this number in the expression of  $P_s$ .

The pseudo-code of this adapted 1/5 rule for tRP-Ens-EDA is given in Algorithm 3. The algorithm first counts the number of improved individuals between the sorted fitness vectors of the two latest generations, and then calculates the success probability  $P_s$ . If  $P_s \leq 1/5$ , then the adaptive method decreases the value of  $c$  (defined as before in eq. (5)) and thus increases  $\nu$  (defined in eq. (6)) for the purpose of exploitation. Otherwise  $c$  is increased, thus  $\nu$  is decreased, for exploration. This method shares with AVS the schemes of updating  $c$  and thresholding  $\nu$ , namely it increases or decreases  $c$  with a coefficient  $\eta \in (0, 1)$ , keeps  $\nu \in [5, 124]$  and resets  $\nu$  to 5 if it goes beyond the range in order to stimulate exploration.

---

**Algorithm 3:** 1/5-success rule based adaptive method.

---

**Input:** Population size  $N$ ; Index of current generation  $t$ ; Degree of freedom  $\nu^t$ , Sorted fitnesses of the latest 2 populations  $\mathbf{f}^t$  &  $\mathbf{f}^{t-1}$ ; Coefficient  $\eta$

**Output:** degree of freedom  $\nu^{t+1}$

```

1  $c \leftarrow 1 + \frac{6}{\nu^{t-4}}$ .
2  $N_p \leftarrow 0$ . //counter of improved individuals.
3 for  $i=1:N$  do
4   if  $f_i^t$  is better than  $f_i^{t-1}$  then
5      $N_p \leftarrow N_p + 1$ .
6  $P_s \leftarrow N_p/N$ . //success probability
7 if  $P_s \leq 1/5$  then
8    $c \leftarrow c\eta$ . //c gets smaller,  $\nu$  gets larger
9 else
10   $c \leftarrow c/\eta$ . //c gets larger,  $\nu$  gets smaller
11  $\nu^{t+1} \leftarrow \frac{6}{c-1} + 4$ .
12 if  $\nu^{t+1} < 5$  or  $\nu^{t+1} > 124$  then
13    $\nu^{t+1} \leftarrow 5$ . //to stimulate exploration
```

---

### 3.3 Adaptive degree of freedom (ADF)

Our third adaptation method will be a scheme where some backtracking is possible. The parameter updating decisions are made after trying out a list of alternative choices and observing their effects on the best fitness within each generation. We drew inspiration from [29] initially, and devise our own adaptation

rules to fit the problem, bearing in mind that our parameter of interest must obey  $\nu \geq 5$ . The pseudo-code of our proposed adaptive method is given in Algorithm 4.

---

**Algorithm 4:** Adaptive degree of freedom (ADF)

---

**Input:**  $\vartheta$  (initial list of  $\nu$ -values in ascending order)

**Output:**  $\nu$

```

1  $L \leftarrow |\vartheta|$ 
2 for  $i = 1 : L$  do
3   Run steps 11-17 and 5 of Algorithm 1 with  $\nu_i$ .
4    $f_i \leftarrow$  best fitness from step 5 of Algorithm 1.
5 if  $\min(\mathbf{f}) == \max(\mathbf{f})$  then
6    $\vartheta_1 \leftarrow \max(5, \text{round}(\vartheta_1/2))$ .
7   for  $j = 2 : L$  do
8      $\vartheta_j \leftarrow \vartheta_j * 2$ .
9 else
10   $\ell \leftarrow \arg \min(\mathbf{f})$ ;  $\nu \leftarrow \vartheta_\ell$ 
11  if  $\ell == 1$  then
12     $\vartheta_1 \leftarrow \max(5, \text{round}(\vartheta_1/2))$ .
13    for  $j = 2 : L$  do
14       $\vartheta_j \leftarrow \max(5, \text{round}((\vartheta_{j-1} + \vartheta_j)/2))$ .
15  else
16    for  $j = 1 : \ell - 1$  do
17       $\vartheta_j \leftarrow \max(5, \text{round}((\vartheta_j + \vartheta_{j+1})/2))$ .
18     $\vartheta_\ell \leftarrow \vartheta_\ell + \max(5, \text{round}(\vartheta_\ell/2))$ .
19    for  $j = \ell + 1 : L$  do
20       $\vartheta_j \leftarrow \max(5, \text{round}((\vartheta_{j-1} + \vartheta_j)/2))$ .
```

---

The adaptive method given in Algorithm 4 maintains a list of  $L$  different alternative values of  $\nu$  in ascending order, and in each generation it creates a separate trial (mock) new-generation using tRP-Ens-EDA with each of the  $L$  different  $\nu$  values on the current list. The list  $L$  is initialised by taking the first  $\nu$  values to be 6, 7, ..., 6 +  $L$  - 1 (all close to the smallest allowable value). The list is then updated based on the current best value, to be tried in the next generation. If the current list of  $\nu$ -values produces no difference to the fitness, then the list of values is spread out more (by decreasing the smallest and increasing all other  $\nu$ -values). Otherwise the algorithm updates the list of values to be tried in the next generation by bringing them closer to the current best  $\nu$ -value.

This adaptation scheme is direct and generic; however it has a tuning parameter that presents a tradeoff: the length of the list of the  $\nu$ -values to try. Obviously, the longer the list, the higher the chance to find a good value within the generation, but the price is a quicker drainage of the budget of fitness evaluations. We experimented with list lengths of  $L = 2$  and  $L = 5$  as representatives of this algorithm in our early work [18], and had found  $L = 2$  to be a good value. Therefore we employed this in the experiments reported here, and will refer to our adaptation scheme as ADF(2df).

It follows from the construction of this method that ADF uses  $L$  times as many fitness evaluations per generation than a fixed choice of  $\nu$  would do, while the two adaptation heuristics presented in the previous subsections keep the same per-generation budget of fitness evaluations – since they make adaptation decisions sequentially, without the possibility of backtracking a parameter value after observing its effect on the fitness. Therefore we expect ADF to be more advantageous in terms of reaching a better fitness value when a sufficiently large budget is available, whereas the sequential heuristics may be preferable when the budget is relatively small. This will be assessed in more detail in the experimental section.

## 4 Experiments

Initial experiments have indicated that the best choice of degree of freedom parameter is problem dependent [18]. Therefore we turn our attention to the adaptation approaches. We conduct a thorough set of experiments to assess our tRP-Ens-EDA with adaptation schemes on a battery of large scale multi-modal benchmark problems, from the 1000-dimensional CEC2010-2012 competition test suite, as described in [31]. All problems are minimisations. Table 1 lists the test functions used in our experiments.

**Table 1.** 1000-dimensional multi-modal test functions from the CEC2010-2012 collection.

Problem	Name
F1	Shifted Rastrigin’s function
F2	Shifted Ackley’s function
F3	Single-group Shifted & m-rotated Rastrigin’s func.
F4	Single-group Shifted & m-rotated Ackley’s func.
F5	$\frac{D}{2B}$ -group Shifted & m-rotated Rastrigin’s func.
F6	$\frac{D}{2B}$ -group Shifted & m-rotated Ackley’s func.
F7	$\frac{D}{2B}$ -group Shifted & m-dimensional Rosenbrock’s func.
F8	$\frac{D}{2B}$ -group Shifted & m-rotated Rastrigin’s func.
F9	$\frac{D}{2B}$ -group Shifted & m-rotated Ackley’s func.
F10	$\frac{D}{2B}$ -group Shifted & m-dimensional Rosenbrock’s func.
F11	Fully nonseparable Rosenbrock

In each experiment we initialise the population randomly in the search space. All the results in comparative experiments are obtained through fresh runs of the methods, except MA-SW-Chains whose results were obtained from the literature.

### 4.1 Results and analysis

#### 4.1.1 Comparative assessment of our adaptation schemes

In the first set of experiments we are interested to find out whether our adaptive schemes can outperform the existing RP-Ens-EDA, and to learn about their respective advantages and disadvantages.

Table 2 provides a statistical analysis of comparisons on the 1000-dimensional benchmarks. For each function, we compare our three different adaptive methods and the existing RP-Ens-EDA of [4] ( $\nu = \infty$ ). We report the average and standard deviation of the best fitness achieved with the evaluation budget of  $6 \times 10^5$  fitness evaluations – this is considered as medium budget in the benchmark problem description.

We perform a two-tailed t-test for each pair of methods involved in this comparison, and the symbols in the rightmost column of the table record statistically significant differences detected at the 0.05 level, as follows: A symbol ‘+’ indicates that the method in the corresponding row performed significantly better than the method in the corresponding column. Likewise, ‘-’ means that the method in the corresponding row per-

formed significantly worse than the method in the corresponding column. Recall that all the test problems are minimisation problems, so lower fitness means better performance. The symbol  $\emptyset$  appears where statistical comparison is not applicable. Positions left blank in the last column mean that no statistically significant difference was detected between the pair of methods tested.

**Table 2.** Statistical comparison of our adaptation schemes and the existing RP-Ens-EDA [4] (that is,  $\nu = \infty$ ) on 1000-dimensional test functions, with a budget of  $6 \times 10^5$  function evaluations.

	Method	Mean	Std	t test vs.		
				AVS	1/5	ADF(2df)
F1	AVS	750.318	57.9958	$\emptyset$		-
	1/5	758.636	48.709		$\emptyset$	-
	ADF(2df)	579.000	24.000	+	+	$\emptyset$
	$\nu = \infty$	784.21	76.017			-
F2	AVS	2.518e-13	4.390e-15	$\emptyset$		
	1/5	2.527e-13	4.390e-15		$\emptyset$	
	ADF(2df)	2.480e-13	4.760e-15			$\emptyset$
	$\nu = \infty$	2.537e-13	4.8104e-15			
F3	AVS	1.308e+07	3.486e+06	$\emptyset$	+	+
	1/5	2.149e+07	9.803e+06	-	$\emptyset$	+
	ADF(2df)	3.040e+08	9.490e+06	-	-	$\emptyset$
	$\nu = \infty$	1.2397e+07	3.1864e+06		+	+
F4	AVS	674.339	126.295	$\emptyset$	+	-
	1/5	1.417e+04	2023.35	-	$\emptyset$	-
	ADF(2df)	548.00	75.00	+	+	$\emptyset$
	$\nu = \infty$	117.78	17.151	+	+	+
F5	AVS	765.787	52.4462	$\emptyset$		-
	1/5	785.701	55.9506		$\emptyset$	-
	ADF(2df)	601.00	26.60	+	+	$\emptyset$
	$\nu = \infty$	832.18	62.543	-	-	-
F6	AVS	34.3466	8.33093	$\emptyset$		
	1/5	35.513	10.6229		$\emptyset$	
	ADF(2df)	35.000	9.490			$\emptyset$
	$\nu = \infty$	41.664	8.7003	-	-	-
F7	AVS	1.3306e+06	7.169e+04	$\emptyset$		+
	1/5	1.3305e+06	6.865e+04		$\emptyset$	+
	ADF(2df)	1.400e+06	5.670e+04	-	-	$\emptyset$
	$\nu = \infty$	1.4442e+06	5.3945e+04	-	-	-
F8	AVS	765.966	46.5231	$\emptyset$		-
	1/5	742.17	45.8631		$\emptyset$	-
	ADF(2df)	599.000	24.600	+	+	$\emptyset$
	$\nu = \infty$	802.11	49.957	-	-	-
F9	AVS	78.5521	14.0886	$\emptyset$		+
	1/5	78.364	13.6203		$\emptyset$	+
	ADF(2df)	91.800	15.700	-	-	$\emptyset$
	$\nu = \infty$	90.481	17.178	-	-	-
F10	AVS	4.363e+04	8810.53	$\emptyset$		
	1/5	4.011e+04	1.106e+04		$\emptyset$	
	ADF(2df)	4.610e+04	9.72e+03			$\emptyset$
	$\nu = \infty$	5.4105e+04	1.0877e+04	-	-	-
F11	AVS	1022.22	63.4622	$\emptyset$		
	1/5	987.983	1.07683		$\emptyset$	+
	ADF(2df)	1.10e+03	67.6		-	$\emptyset$
	$\nu = \infty$	1614.7	249.44	-	-	-

From these results we observe the following.

- In 9 (out of 11) of the test problems, no statistically significant difference is detected between the AVS, and the 1/5 adaptation heuristics.

- In the remaining 2 problems, AVS is statistically superior to 1/5. On further investigation we observed that the 1/5 approach has periods of stagnation where AVS is able to progress. This is most likely because the definition of the 1/5 in the context of EDA is rather ad hoc.
- AVS wins with statistical significance against the existing RP-Ens-EDA in 7 (out of 11) test problems, and only loses on 1 (F4). The remaining 3 cases are ties.
- ADF performs on-par with AVS at this budget size: 4 wins and 3 losses against AVS, and it outperforms the existing RP-Ens-EDA in 7 of the functions; it is outperformed by RP-Ens-EDA in 2 functions (F4 and F3), 2 ties.

We then ran the two most successful adaptation methods, ADF(2df) and AVS, versus the existing RP-Ens-EDA for a larger budget size of  $3 \cdot 10^6$ . Fig. 1 shows the detailed trajectories, in order to visualise the behaviour of these methods and gain insight into the reasons for the observed differences. In Fig. 1 the best fitness is plotted against the elapsed number of function evaluations, averaged over 25 independent restarts. For clarity of visual inspection, these are zoomed in to the interesting ranges of the axes: the fitness values of the first few generations are not shown, and the plotting halts when no more interesting change occurs.

We see from the trajectory plots on figure 1 that, given a sufficient budget, the direct adaptation method ADF(2df) often reaches a better fitness value with a slight delay. This makes good sense intuitively too, based on the construction of the algorithm that spends part of its resources on exploring. Interestingly, it often goes through a period of stagnation before it wins over. For instance, on F3, with the medium budget size reported in Table 2 it was outperformed by all

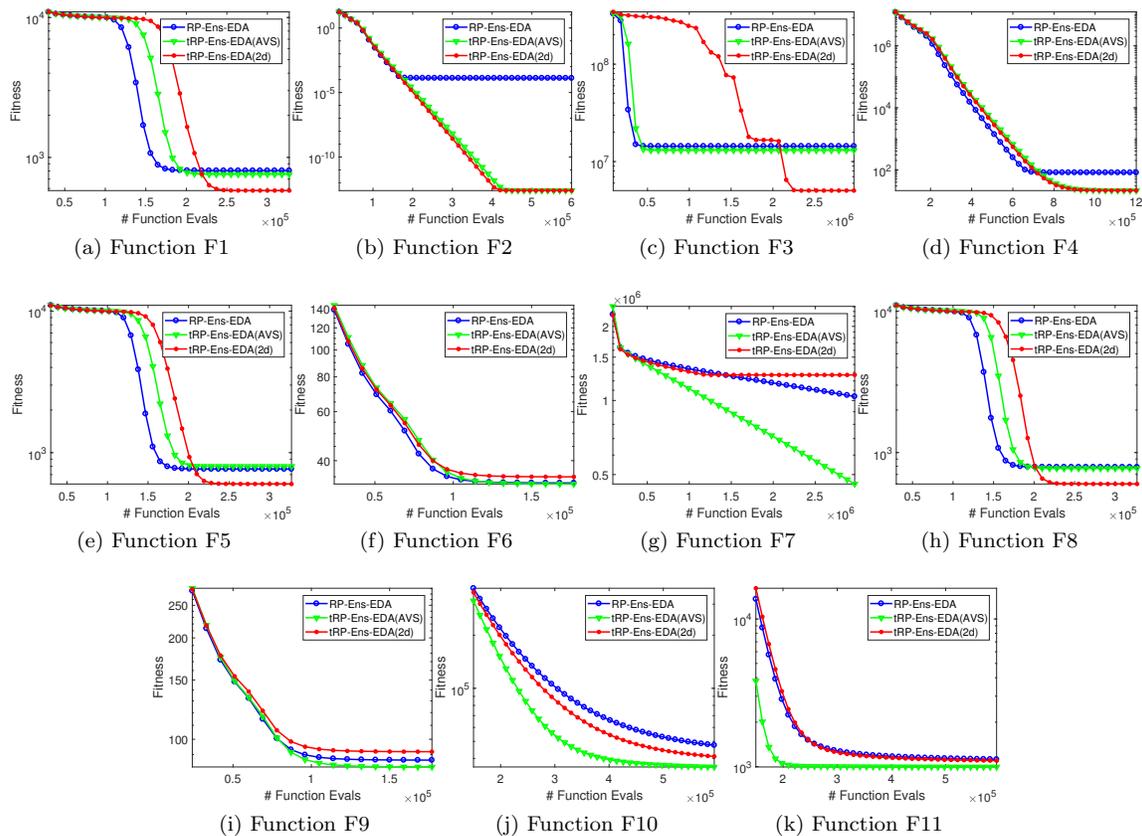


Fig.1. Fitness trajectories of the RP-Ens-EDA, tRP-Ens-EDA(AVS) and tRP-Ens-EDA(2df) for a maximum of  $3 \times 10^6$  function evaluations. For a statistical analysis of significance of the differences observed halfway through the budget after  $6 \times 10^5$  function evaluations, see Table 2.

other methods, including the previously existing RP-Ens-EDA, whereas from the trajectory plot we see that it does win eventually, given a larger budget, whereas the Gaussian variant makes no further progress. Indeed, we see from the trajectories in Fig. 1 that the Gaussian RP-Ens-EDA is more prone to premature convergence than our adaptive tRP-Ens-EDA variants. More differences between ADF and AVS will be seen in multi-comparison tests in the presence of other competing state-of-the-art algorithms in section 4.1.2.

#### 4.1.2 Comparison with state of the art methods on 1000 dimensional problems

We compare our best performing adaptive degree of freedom methods tRP-Ens-EDA(ADF) and tRP-Ens-EDA(AVS) with the existing state-of-the-art on the

1000 dimensional CEC 2010-2012 competition multi-modal benchmark problems.

Tables 3-4 and 5-6 summarise the results from experiments with two different budgeted sizes respectively, from 25 repeated independent runs each. Average and standard deviation of the best fitness values achieved are reported. By the Friedman’s test and a subsequent multi-comparison, the means and standard deviations of the best performing algorithms are marked in bold. The two budget sizes tested are  $6 \cdot 10^5$  and  $3 \cdot 10^6$  – that is, a medium size budget and a large budget, according to the definitions in the benchmark suite.

The results are summarised in Tables 3- 4 and 5- 6. We see from these tables, and also from figure 1, that our tRP-Ens-EDA(ADF) is most competitive in the high

budget setting. In particular,

- In the medium budget setting ( $0.6 \cdot 10^6$  function evaluations), tRP-Ens-EDA(ADF) wins with statistical significance over all competitors tested in two of the test functions in the medium budget setting. AVS wins over all competitors tested in one function. In fact, tRP-Ens-EDA(AVS) does better than tRP-Ens-EDA(ADF) in terms of achieving lower fitness on a larger number of functions, but statistically significant out-performance is not observed in most of these cases in the multi-comparison test.
- In the large budget setting ( $3 \cdot 10^6$  function evaluations), tRP-Ens-EDA(ADF) wins over all other methods tested – including the winner of the CEC’10-’12 competitions [15] – in 4 of the test functions (out of 11 tested).

We can therefore conclude that a direct adaptation scheme like ADF turns most beneficial when we allow slightly larger budget sizes. In medium or limited budget settings we can recommend the AVS method.

Naturally, one cannot expect any method to perform best on all problems in the light of well known no free lunch theorems. However, the observed empirical results together with their statistical significance analysis support the conclusion that the methods presented in this paper are competitive with state of the art large-scale optimisation heuristics, and therefore make a worthwhile addition to the practitioner’s toolbox for use in difficult high dimensional problems that no specialised algorithm is known to be able to solve. Moreover, our algorithm construction comes with a good intuition of its working strategy, which consists in a controlled increase of exploration in a similar spirit as the adaptive variance scaling idea previously pursued in [27], but

now supporting much larger scale problems through an ensemble of RPs.

#### 4.1.3 Scalability of ADF(2df)

Since our ADF(2df) is the most expensive in terms of function evaluations per generation, our final set of experiments is to assess its scalability. We measure the number of function evaluations needed (search cost) to reach a specified distance from the global optimum as the problem dimension varies. We fix the value to reach (VTR) to  $10^{-3}$ , and vary the dimensionality of the problem  $d \in [50, 1000]$ . We repeat the same experiment with three other choices of VTR:  $10^{-2}$ ,  $10^2$  and  $10^3$  to avoid biasing the conclusions towards a particular choice of the VTR. Throughout, the budget was fixed to  $3 \cdot 10^6$  so the algorithm stops either when the budget is exhausted or upon reaching the VTR. All other parameters are kept at their default values.

Figure 2 displays the average number of function evaluations as computed from the successful runs out of 25 independent repetitions for each problem, for each dimension tested. That is, runs terminated by exhausting the budget do not contribute to these plots.

From figure 2 we observe a linear shape on the scalability measurements (dashed lines on the log-log plots). The black dotted line corresponds to the slope of 1 (that is, linear scaling), for reference. This means that our tRP-Ens-EDA with ADF scheme scales at most polynomially (with degree indicated by the slope of lines) in the dimension of the problem, which matches the best scaling known for sep-CMA-ES [7] and the original RP-Ens-EDA [4].

## 5 Conclusions and future work

We devised new adaptive approaches for high dimensional continuous black-box optimisation which extend and generalise the random projection ensemble based

**Table 3.** Mean of best fitness values from 25 independent runs, in comparison with state of the art under equal budget of  $0.6 \cdot 10^6$  function evaluations.

	MA-SW-Chains	CCVIL	MLCC	sep-CMA-ES	EDA-MCC	tRP-Ens-EDA(AVS)	RP-Ens-EDA	tRP-Ens-EDA(ADF)
F1	2.67E+03	<b>1.77E+01</b>	6.22E+03	8.15E+03	1.24E+03	7.57E+02	8.05E+02	5.79E+02
F2	3.84E+00	1.77E+01	1.45E+01	2.15E+01	3.86E+00	2.50E-13	1.30E-04	2.48E-13
F3	2.17E+08	9.79E+08	4.19E+08	1.41E+08	1.84E+07	1.30E+07	1.44E+07	3.04E+08
F4	8.14E+04	2.12E+07	1.41E+07	1.04E+07	<b>1.78E+01</b>	6.69E+02	2.51E+02	5.48E+02
F5	3.22E+03	1.31E+04	1.10E+04	8.36E+03	1.38E+03	8.01E+02	7.66E+02	<b>6.01E+02</b>
F6	3.83E+01	2.32E+02	2.16E+02	2.25E+02	3.42E+01	2.30E+01	3.33E+01	3.50E+01
F7	<b>4.34E+03</b>	1.84E+11	1.21E+09	5.02E+05	1.89E+08	1.33E+06	1.43E+06	1.40E+06
F8	3.19E+03	1.86E+04	1.84E+04	8.13E+03	1.49E+03	7.72E+02	7.89E+02	<b>5.99E+02</b>
F9	1.02E+02	4.26E+02	4.14E+02	4.30E+02	7.92E+01	8.27E+01	8.67E+01	9.18E+01
F10	<b>5.53E+03</b>	4.76E+11	2.67E+10	2.19E+04	1.65E+09	3.58E+04	4.74E+04	4.61E+04
F11	1.21E+03	4.48E+11	5.26E+10	1.16E+03	2.16E+09	<b>9.94E+02</b>	1.12E+03	1.10E+03

**Table 4.** Standard deviation of best fitness values from 25 independent runs, in comparison with state of the art, under equal budget of  $0.6 \cdot 10^6$  function evaluations.

	MA-SW-Chains	CCVIL	MLCC	sep-CMA-ES	EDA-MCC	tRP-Ens-EDA(AVS)	RP-Ens-EDA	tRP-Ens-EDA(ADF)
F1	1.63E+02	<b>1.40E+01</b>	1.75E+03	3.63E+02	4.77E+01	6.76E+01	7.17E+01	2.40E+01
F2	2.13E-01	1.97E+00	1.90E+00	2.07E-02	3.66E-01	6.39E-015	2.24E-06	4.76E-15
F3	8.56E+07	7.60E+07	1.76E+08	3.00E+07	3.39E+06	4.09E+06	3.29E+06	9.49E+06
F4	2.84E+05	3.76E+05	8.29E+06	2.33E+06	<b>1.90E-01</b>	1.30E+02	2.68E+01	7.50E+01
F5	1.85E+02	3.29E+02	1.97E+03	3.85E+02	6.28E+01	9.17E+01	5.91E+01	<b>2.66E+01</b>
F6	7.23E+00	8.08E-01	1.79E+01	5.05E+00	3.12E+00	9.48E+00	7.31E+00	9.49E+00
F7	<b>3.21E+03</b>	1.23E+11	2.15E+09	8.66E+04	7.86E+07	5.45E+04	4.72E+04	5.67E+04
F8	1.46E+02	3.76E+02	1.02E+04	3.49E+02	7.53E+01	3.79E+01	4.85E+01	<b>2.46E+01</b>
F9	1.42E+01	9.41E-01	1.22E+01	2.88E+00	8.38E+00	1.29E+01	1.43E+01	1.57E+01
F10	<b>3.94E+03</b>	2.86E+11	5.31E+10	1.12E+04	3.97E+08	5.43E+03	1.12E+04	9.72E+03
F11	1.42E+02	2.64E+11	7.80E+10	1.34E+02	4.97E+08	<b>1.55E+01</b>	1.04E+02	6.76E+01

EDA by means of using unconventional random matrices. We have analytically shown that this results in enlarging the covariance of the high dimensional search distribution, which then increases exploration. Hence our adaptive schemes in fact implement novel means of covariance adaptation. We have demonstrated superior performance against the original version RP-Ens-EDA on both large and medium budget instances, and with a state of the art performance in comparison with a number of existing methods on 1000 dimensional problems. Future work remains to investigate other aggregation schemes, and to extend the method to multi-objective problems. A very interesting question for future work would be to analyse the dynamics and the running time of such high dimensional EDA optimisation algorithms.

## Acknowledgments

The authors would like to thank Qi Xu (MSc student at Birmingham) for implementing the AVS method and carrying out part of the experiments for Tables 5 and 6. The work of MLS was partly funded by a PhD scholarship from the Islamic Development Bank. The work of AK is funded by Engineering and Physical Sciences Research Council under Fellowship Grant EP/P004245/1.

## References

- [1] P. Larrañaga, J.A. Lozano. *Estimation of Distribution Algorithms: A new tool for Evolutionary Computation*. Norwell, MA: Kluwer Academic Publishers, 2001.
- [2] B. Yuan, M. Gallagher. On the importance of di-

**Table 5.** Mean of best fitness values from 25 independent runs, in comparison with state of the art, under equal budget of  $3 \cdot 10^6$  function evaluations.

	MA-SW-Chains	CCVIL	MLCC	sep-CMA-ES	EDA-MCC	tRP-Ens-EDA(AVS)	RP-Ens-EDA	tRP-Ens-EDA(ADF)
F1	<b>2.10E-14</b>	4.00E-07	5.57E-01	5.68E+03	1.24E+03	7.56E+02	8.05E+02	5.94E+02
F2	7.28E-13	4.00E-07	9.88E-13	2.15E+01	3.86E+00	2.52E-13	1.30E-04	<b>2.46E-13</b>
F3	1.68E+08	5.51E+08	3.84E+08	1.19E+08	1.84E+07	1.30E+07	1.44E+07	<b>4.99E+06</b>
F4	8.14E+04	4.14E+05	1.62E+07	6.39E+06	<b>1.78E+01</b>	2.13E+01	8.05E+01	3.20E+01
F5	2.07E+03	1.43E+03	3.43E+03	6.28E+03	1.38E+03	8.01E+02	7.66E+02	<b>5.96E+02</b>
F6	3.80E+01	<b>7.44E+00</b>	1.98E+02	2.12E+02	3.42E+01	3.30E+01	3.33E+01	3.10E+01
F7	1.25E+03	2.98E+11	2.08E+03	<b>2.94E+02</b>	1.89E+08	4.51E+05	1.03E+06	1.00E+06
F8	2.74E+03	2.78E+03	7.11E+03	6.76E+03	1.49E+03	7.72E+02	7.89E+02	<b>6.05E+02</b>
F9	9.98E+01	<b>1.31E+01</b>	3.76E+02	4.21E+02	7.92E+01	8.27E+01	8.67E+01	9.18E+01
F10	1.30E+03	6.42E+11	7.09E+03	<b>9.16E+02</b>	1.65E+09	2.87E+04	4.02E+04	4.05E+04
F11	1.07E+03	1.75E+11	2.05E+03	9.04E+02	2.16E+09	9.89E+02	1.04E+03	1.06E+03

**Table 6.** Standard deviation of best fitness values from 25 independent runs, in comparison with state of the art, under equal budget of  $3 \cdot 10^6$  function evaluations.

	MA-SW-Chains	CCVIL	MLCC	sep-CMA-ES	EDA-MCC	tRP-Ens-EDA(AVS)	RP-Ens-EDA	tRP-Ens-EDA(ADF)
F1	<b>1.99E-14</b>	0.00E+00	2.21E+00	4.89E+02	4.77E+01	6.76E+01	7.17E+01	1.76E+01
F2	3.40E-13	6.32E-07	3.70E-12	1.08E-01	3.66E-01	6.30E-15	2.33E-06	<b>4.58E-15</b>
F3	1.04E+08	1.55E+08	6.93E+07	2.92E+07	3.39E+06	4.09E+06	3.29E+06	<b>1.32E+06</b>
F4	2.84E+05	6.54E+05	4.97E+06	3.85E+06	<b>1.90E-01</b>	2.82E-02	2.13E+00	5.32E-01
F5	1.44E+02	6.34E+01	8.72E+02	2.51E+02	6.28E+01	9.17E+01	5.91E+01	<b>1.71E+01</b>
F6	7.35E+00	<b>2.41E+00</b>	6.98E-01	6.16E+00	3.12E+00	9.48E+00	7.31E+00	1.09E+01
F7	5.72E+02	8.57E+11	7.27E+02	<b>9.20E+01</b>	7.86E+07	6.82E+04	5.74E+04	5.74E+04
F8	1.22E+02	8.78E+01	1.34E+03	2.76E+02	7.52E+01	3.79E+01	4.84E+01	<b>2.74E+01</b>
F9	1.40E+01	<b>2.92E+00</b>	4.71E+01	1.59E+01	8.39E+00	1.29E+01	1.43E+01	1.57E+01
F10	4.36E+02	2.19E+12	4.77E+03	<b>1.94E+02</b>	3.96E+08	5.13E+03	9.86E+03	9.82E+03
F11	7.29E+01	8.77E+11	1.80E+02	3.91E+01	4.97E+08	1.06E+01	5.08E+01	5.45E+01

versity maintenance in estimation of distribution algorithms. In *Proceedings of the Genetic and Evolutionary Computation Conference, June 2005*, pp. 719-729.

[3] J. H Friedman. *An overview of predictive learning and function approximation*. Springer, 1994, Vol. 136.

[4] A. Kabán, J. Bootkrajang, R. J. Durrant. Toward large-scale continuous EDA: A random matrix theory perspective. *Evolutionary Computation, 2016*, 24(2): 255-291.

[5] W. Dong, X. Yao. Covariance matrix repairing in gaussian based edas. In *IEEE Congress on Evolutionary Computation, September 2007*, pp. 415-422.

[6] T.K, H. Iba. Linear and combinatorial optimizations by estimation of distribution algorithms. In *Proceedings of the 9th MPS Symposium on Evolutionary Computation, July 2002*, pp. 99-106.

[7] R. Ros, N. Hansen. A simple modification in cma-es achieving linear time and space complexity. In *Parallel Problem Solving from Nature-PPSN X, September 2008*, pp. 296-305.

[8] P. Bosman. On empirical memory design, faster selection of Bayesian factorizations and parameter-free Gaussian EDAs. In *Proceedings of Genetic and Evolutionary Computation Conference, July 2009*, pp. 389-396.

[9] J.S. De Bonet, C.L. Isbell, Jr. P. Viola. Mimic: Finding optima by estimating probability densi-

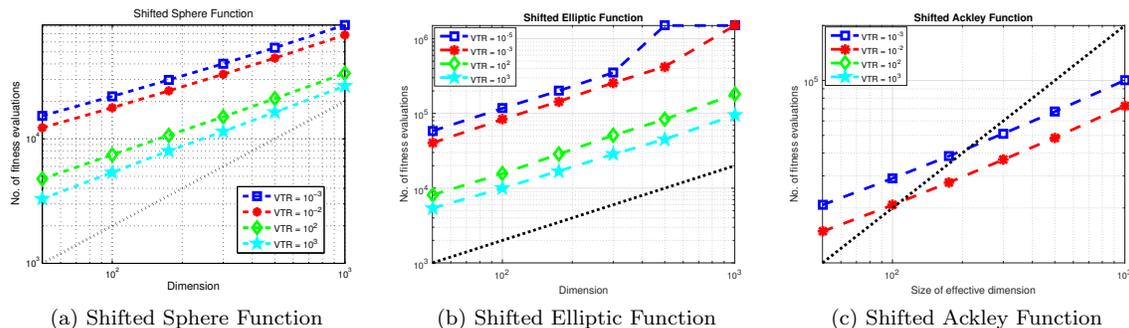


Fig.2. Number of function evaluations taken by successful runs of our tRP-Ens-EDA with ADF to reach a pre-specified value to reach (VTR) as the problem dimensionality is varied in  $d \in [50, 1000]$ . The markers represent averages computed from 25 independent repetitions. The black dotted line corresponds to linear scaling (slope = 1).

ties. In *Neural Information Processing Systems*. 1997, 9(1): 424-430.

[10] P.A.N. Bosman, D. Thierens. An algorithmic framework for density estimation based evolutionary algorithms. Technical report, 1999, UU-CS.

[11] R. Armañanzas, I. Inza, R. Santana, Y. Saeys, J.L. Flores, J.A. Lozano, Y. Van de Peer, R. Blanco, V. Robles, C. Bielza, P. Larrañaga. A review of estimation of distribution algorithms in bioinformatics. *BioData Mining*, 2008, 1(6).

[12] K. Weicker, N. Weicker. On the improvement of co-evolutionary optimizers by learning variable inter-dependencies. In *Proceedings of the Congress on Evolutionary Computation, July 1999*, pp.1627-1632.

[13] W. Dong, T. Chen, P. Tiño, X. Yao. Scaling up estimation of distribution algorithm for continuous optimisation. *IEEE Transaction of Evolutionary Computation*. 2013, 17(6): 797 - 822.

[14] Z. Yang, K. Tang, X. Yao. Multilevel cooperative co-evolution for large scale optimization. In *IEEE World Congress on Computational Intelligence, June 2008*, pp. 1663 - 1670.

[15] D. Molina, M. Lozano, A. M. Sánchez, F. Herrera. MA-SSW-Chains: Memetic algorithm based on local search chains for large scale continuous global optimization. *Soft Computing*, 2011, 15(11): 2201-2220.

[16] S. Y. Shin, D.Y. Cho, B. Tak Zhang. Function optimization with latent variable models, in a. ochoa et al. (eds.). In *Proceedings of the Third International Symposium on Adaptive Systems ISAS, June 2001*, pp. 145-152.

[17] M. L. Sanyang, A. Kabán. Remeda: Random embedding EDA for optimising functions with intrinsic dimension. In *Parallel Problem Solving from Nature - PPSN XIV*, September 2016, pp. 859-868.

[18] M. L. Sanyang, A. Kabán. Heavy tails with parameter adaptation in random projection based continuous EDA. In *Proceedings of IEEE Congress on Evolutionary Computation, May 2015*, pp. 2074-2081.

[19] X. Yao, Y. Liu, G. Lin. Evolutionary programming made faster. *IEEE Transaction on Evolutionary Computation*, 1999, 3(2): 82-102.

[20] B. Gao, I. Wood. Multivariate t distribution, archive and mutation based estimation of distri-

- bution algorithm. *ANZIAM Journal, [S.l.], 2014, 54(SUPPL): pp. 720-746.*
- [21] B. Gao. Estimation of distribution algorithms for single- and multi-objective optimization. *PhD Thesis, University of Queensland,, 2014.*
- [22] M.L. Sanyang, R.J. Durrant, A. Kabán. How effective is cauchy-eda in high dimensions? In *IEEE Congress on Evolutionary Computation, CEC Vancouver, BC, Canada, July 2016*, pp. 3409-3416.
- [23] Duc-Cuong Dang, Per Kristian Lehre, Phan Trung Hai Nguyen. Level-based analysis of the univariate marginal distribution algorithm. *Algorithmica, 2019, 81(2): 668-702.*
- [24] A. Kabán. On compressive ensemble induced regularisation: How close is the finite ensemble precision matrix to the infinite ensemble? In *International Conference on Algorithmic Learning Theory, October 2017*, pp. 617-628.
- [25] A. Kabán. New bounds for compressive linear least squares regression. In *The 17-th International Conference on Artificial Intelligence and Statistics, April 2014*, pp. 448-456.
- [26] D. Achlioptas. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *Journal of Computer and System Sciences, 2003, 66(4): 671-687.*
- [27] J. Grahl, P. A. Bosman, F. Rothlauf. The correlation-triggered adaptive variance scaling idea. In *Proceedings of the annual conference on Genetic and evolutionary computation, July 2006*, pp. 397-404.
- [28] A.E. Eiben, R. Hinterding, Z. Michalewicz. Parameter control in evolutionary algorithms. *IEEE Transaction on Evolutionary Computation, 1999, 3(2): 124-141.*
- [29] Y.Y. Wong, K.H. Lee, K.S. Leung, C.W. Ho. A novel approach in parameter adaptation and diversity maintenance for genetic algorithms. *Soft Computing, 2003, 7(8): 506-515.*
- [30] H. G. Beyer, H. P. Schwefel. Evolution strategies- a comprehensive introduction. *Natural computing, 1999, 1(1): 3-52.*
- [31] K. Tang, Xiaodong Li, P. N. Suganthan, Z. Yang, T. Weise. Benchmark functions for the CEC'2010 special session and competition on large scale global optimization. Technical report, Nature Inspired Computation and Applications Laboratory, USTC, China, 2012.
- [32] M. Abramowitz, I. A. Stegun (Eds.). Handbook of mathematical functions with formulas, graphs and mathematical tables (applied mathematics series 55). Washington: NBS, 1964.

## Appendix A

*Proof* [of Proposition 1] Recall from step 14 of Algorithm 1 that the set of projected points in the  $i$ -th subspace is:

$$\mathcal{Y}^{\mathbf{R}_i} = \{\mathbf{R}_i(\mathbf{x}_1 - \boldsymbol{\mu}), \mathbf{R}_i(\mathbf{x}_2 - \boldsymbol{\mu}), \dots, \mathbf{R}_i(\mathbf{x}_N - \boldsymbol{\mu})\}.$$

Conditional on  $\mathbf{R}_i$ , the sample covariance matrix of this is:

$$\boldsymbol{\Sigma}^{\mathbf{R}_i} = \frac{1}{N} \sum_{n=1}^N \mathbf{R}_i(\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{R}_i(\mathbf{x}_n - \boldsymbol{\mu}))^T = \mathbf{R}_i \boldsymbol{\Sigma} \mathbf{R}_i^T.$$

So the samples in step 16 are  $\mathbf{y}_1^{\mathbf{R}_i}, \dots, \mathbf{y}_N^{\mathbf{R}_i} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{\mathbf{R}_i})$ .

To find the distribution of the individuals in  $\mathcal{P}^{new}$  at step 17, we look at the  $n$ -th individual:

$$\mathbf{x}_n := \sqrt{\frac{dM}{k}} \left[ \frac{1}{M} \sum_{i=1}^M \mathbf{R}_i^T \mathbf{y}_n^{\mathbf{R}_i} \right] + \boldsymbol{\mu}. \quad (7)$$

Conditionally on  $\mathbf{R}_i, i = 1 : M$ , this is a linear combination of independent Gaussian random variables, which is again a Gaussian<sup>†</sup>. Hence,  $\mathbf{x}_n$  is Gaussian distributed, with mean  $\boldsymbol{\mu}$  (since  $\mathbf{y}_n^{\mathbf{R}_i}$  has zero mean), and we compute its covariance below.

In equation (7), denote  $\mathbf{A}_i := \sqrt{\frac{dM}{k}} \frac{1}{M} \mathbf{R}_i^T$ , then we have that  $\mathbf{y}_n^{\mathbf{R}_i} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_i \boldsymbol{\Sigma} \mathbf{R}_i^T)$ . So,

$$\mathbf{A}_i \mathbf{y}_n^{\mathbf{R}_i} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_i \mathbf{R}_i \boldsymbol{\Sigma} \mathbf{R}_i^T \mathbf{A}_i^T) \quad (8)$$

for all  $n = 1, \dots, N$ . Replacing  $\mathbf{A}_i$  in (8), we have the  $i$ -th summand in eq. (7) is:

$$\mathbf{A} \mathbf{y}_n^{\mathbf{R}_i} \sim \mathcal{N}\left(\mathbf{0}, \sqrt{\frac{dM}{k}} \frac{1}{M} \mathbf{R}_i^T \mathbf{R}_i \boldsymbol{\Sigma} \mathbf{R}_i^T \sqrt{\frac{dM}{k}} \frac{1}{M} \mathbf{R}_i\right),$$

which simplifies to

$$\mathcal{N}\left(\mathbf{0}, \frac{d}{kM} \mathbf{R}_i^T \mathbf{R}_i \boldsymbol{\Sigma} \mathbf{R}_i^T \mathbf{R}_i\right).$$

Finally, the summation operation yields the ensemble covariance in the  $d$ -dimensional search space as stated in eq. (1).  $\square$

## Appendix B

Here we give an example to demonstrate that random matrices with negative excess kurtosis can cancel the regularisation effect of the ensemble in some directions of the search space.

Let  $R_{ij} \sim \{-1, +1\}$  with a probability of 1/2 each. Hence we have  $E[R_{ij}^4] = 1$ ,  $E[R_{ij}^2] = 1$ , and therefore the excess kurtosis is  $K = 1 - 3 = -2$ .

For the sake of this example, suppose that  $\text{rank}(\boldsymbol{\Sigma}) = 2$ , so  $\lambda_1 \geq \lambda_2 > 0, \lambda_3 = \dots \lambda_d = 0$ , and suppose the first two eigenvectors of  $\boldsymbol{\Sigma}$  are

$$\begin{aligned} \mathbf{u}_1 &= \left(1/\sqrt{2}, 1/\sqrt{2}, 0, \dots, 0\right), \\ \mathbf{u}_2 &= \left(-1/\sqrt{2}, 1/\sqrt{2}, 0, \dots, 0\right), \end{aligned}$$

and the remaining eigenvectors have their first two coordinates 0. This is well defined, since one can easily verify that  $\mathbf{u}_1 \perp \mathbf{u}_2$ ,  $\|\mathbf{u}_1\| = \|\mathbf{u}_2\| = 1$ . Denote by  $\mathbf{U}$  the matrix having  $\mathbf{u}_i$  in its rows, with entries  $U_{ai}$ , where  $a, i = 1, \dots, d$ . Now, we have:

$$\begin{aligned} \mathbf{A}_1 &= \begin{bmatrix} \sum_{a=1}^d U_{a1}^4 & 0 & 0 & \dots & 0 \\ 0 & \sum_{a=1}^d U_{a1}^2 U_{a2}^2 & 0 & \dots & 0 \\ 0 & 0 & \sum_{a=1}^d U_{a1}^2 U_{a3}^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sum_{a=1}^d U_{a1}^2 U_{ad}^2 \end{bmatrix} \\ &= \begin{bmatrix} 1/2 & 0 & 0 & \dots & 0 \\ 0 & 1/2 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} = \mathbf{A}_2. \end{aligned}$$

Then, for this example we get:

$$\begin{aligned} K \sum_{i=1}^d \lambda_i \mathbf{A}_i &= -2(\lambda_1 \mathbf{A}_1 + \lambda_2 \mathbf{A}_2) \\ &= -\text{Tr}(\boldsymbol{\Sigma}) \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}. \end{aligned}$$

This cancels the effect of the term  $\text{Tr}(\boldsymbol{\Sigma}) \mathbf{I}_d$  in the first two coordinate directions, leaving  $E_{\mathbf{R}}[\boldsymbol{\Sigma}_{rp}]$  to be a multiple of maximum likelihood covariance estimates in those directions.

## Appendix C

*Proof* [Proof of Proposition 2] Let  $k \in \{2, 4\}$  and  $x \sim t(0, 1, \nu)$  a t-distributed random variable, then by definition,

$$E[x^k] = \int_{-\infty}^{\infty} x^k f(x) dx$$

where  $f(x)$  is the pdf of the t distribution,

$$f(x) = c \left(1 + \frac{x^2}{\nu}\right)^{-\frac{1}{2}(\nu+1)}, \quad (9)$$

<sup>†</sup>Assume  $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \boldsymbol{\Sigma}_x)$  and  $\mathbf{y} \sim \mathcal{N}(\mathbf{m}_y, \boldsymbol{\Sigma}_y)$ , then

$$\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c} \sim \mathcal{N}(\mathbf{A}\mathbf{m}_x + \mathbf{B}\mathbf{m}_y + \mathbf{c}, \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_y\mathbf{B}^T)$$

where

$$c = \frac{1}{\sqrt{\nu}B(\frac{\nu}{2}, \frac{1}{2})} = \frac{\Gamma(\nu/2 + 1/2)}{\sqrt{\nu} \Gamma(\nu/2)\Gamma(1/2)} \quad (10)$$

and  $B(\cdot, \cdot)$  is the beta function [32]. Observing that  $f(x) = f(-x)$ , we can re-write the integral as:

$$E[x^k] = 2 \int_0^\infty x^k f(x) dx. \quad (11)$$

Writing out the expression of  $f$  as defined in eq. (9), we have that eq. (11) equals

$$E[x^k] = 2c \int_0^\infty x^k \left(1 + \frac{x^2}{\nu}\right)^{-\frac{1}{2}(\nu+1)} dx$$

Now, we make the change of variable:  $t = \frac{x^2}{\nu}$ . Hence  $x = (\nu t)^{\frac{1}{2}}$ , and we get:

$$\begin{aligned} E[x^4] &= c\nu^{\frac{k+1}{2}} \int_0^\infty t^{\frac{k+1}{2}-1} (1+t)^{-\frac{1}{2}-\frac{\nu}{2}} dt \\ &= c\nu^{\frac{k+1}{2}} \int_0^\infty t^{\frac{k+1}{2}-1} (1+t)^{-(\frac{k+1}{2})-(\frac{\nu-k}{2})} dt, \end{aligned}$$

where we added and subtracted  $k/2$  in the exponent of the last term so that the integral represents a beta function [32]:

$$B(m+1, n+1) = \int_0^\infty u^m (1+u)^{-(m+n)-2}$$

with  $m := \frac{k}{2} - \frac{1}{2}$  and  $n := \frac{\nu}{2} - 3$ . Therefore we can write:

$$\begin{aligned} E[x^k] &= c\nu^{\frac{k+1}{2}} B\left(\frac{k+1}{2}, \frac{\nu-k}{2}\right) \\ &= c\nu^{\frac{k+1}{2}} \frac{\Gamma(\frac{k+1}{2})\Gamma(\frac{\nu-k}{2})}{\Gamma(\frac{1+\nu}{2})}. \end{aligned}$$

We need to evaluate this for  $k = 4$  and  $k = 2$ . Replacing the expression of  $c$  from eq. (10), we get for  $k = 4$  the following:

$$E[x^4] = \nu^2 \frac{\Gamma(\frac{5}{2})\Gamma(\frac{\nu-4}{2})}{\Gamma(\frac{\nu}{2})\Gamma(\frac{1}{2})},$$

which is finite provided that  $\nu-4 > 0$  (since  $\Gamma(0) = \infty$ ). To simplify, we use the identities  $\Gamma(\frac{5}{2}) = \frac{3}{4}\sqrt{\pi}$  and  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$  [32]. So,

$$E[x^4] = \frac{3\nu^2}{4} \frac{\Gamma(\frac{\nu-4}{2})}{\Gamma(\frac{\nu}{2})}.$$

Furthermore, by a property of the Gamma function,  $\Gamma(x) = (x-1)\Gamma(x-1)$  [32], we have after substitution and simplification:

$$E[x^4] = \frac{3\nu^2}{(\nu-2)(\nu-4)}.$$

Analogously, for  $k = 2$ , we arrive at

$$E[x^2] = \frac{\nu}{\nu-2}.$$

Therefore, the excess kurtosis is:

$$K = \frac{E[x^4]}{(E[x^2])^2} - 3 = \frac{6}{\nu-4}.$$

□

*Proof* [of Corollary 1] Let  $c > 0$  be a constant. Then  $c \cdot x$  has variance  $c^2 \text{var}(x)$ . The excess kurtosis of  $c \cdot x$  is

$$\frac{E[(c \cdot x)^4]}{(E[(c \cdot x)^2])^2} - 3.$$

Taking the constant out, we have

$$\frac{c^4 E[x^4]}{c^4 E[x^2]^2} - 3 = \frac{E[x^4]}{E[x^2]^2} - 3,$$

thus the excess kurtosis did not change. □

**Momodou L. Sanyang** received his Ph.D degree in Computer Science from the University of Birmingham, UK, in 2017. Prior to that, he received his B.Sc. degree in Mathematics and Physics from



The University of The Gambia (UTG), The Gambia, in 2006, and his M.Sc. degree in Information Systems and Applications from the National Tsing Hua University, Hsinchu, in 2010. He is currently a senior lecturer in Computer Science at the University of The Gambia (UTG), and a honorary Research Fellow at the University of Birmingham, UK. His research interests include machine learning and data mining, dimensionality reduction, random projections, evolutionary computation and black box optimisation in high dimensional settings.



**Ata Kabán** is Professor in Computer Science at the University of Birmingham, UK. Her main research interests are in statistical machine learning, data mining, and black-box optimisation in high dimensional settings. She holds a Ph.D in Computer Science (2002) from the University of Paisley,

Scotland, and a Ph.D in Musicology (2000) from the Music Academy "Gh. Dima" of Cluj-Napoca, Romania. She is member of the IEEE CIS Technical Committee on Data Mining and Big Data Analytics, and a vice-chair of the IEEE CIS Task Force on High Dimensional Data Mining. She has been awarded a 5-years EPSRC Fellowship (2017-2022), and a part-time Turing Fellowship (2019-2021).