Natural Gradient Approach for Linearly Constrained Continuous Optimization

Youhei Akimoto¹ and Shinichi Shirakawa²

¹ Faculty of Engineering, Shinshu University, Nagano, Nagano, Japan y_akimoto@shinshu-u.ac.jp
 ² College of Science and Engineering, Aoyama Gakuin University, Sagamihara, Kanagawa, Japan shirakawa@it.aoyama.ac.jp

Abstract. When a feasible set of an optimization problem is a proper subset of a multidimensional real space and the optimum of the problem is located on or near the boundary of the feasible set, most evolutionary algorithms require a constraint handling machinery to generate better candidate solutions in the feasible set. However, some standard constraint handling such as a resampling strategy affects the distribution of the candidate solutions; the distribution is truncated into the feasible set. Then, the statistical meaning of the update of the distribution parameters will change. To construct the parameter update rule for the covariance matrix adaptation evolution strategy from the same principle as unconstrained cases, namely the natural gradient principle, we derive the natural gradient of the log-likelihood of the Gaussian distribution truncated into a linearly constrained feasible set. We analyze the novel parameter update on a minimization of a spherical function with a linear constraint.

1 Introduction

The covariance matrix adaptation evolution strategy (CMA-ES) is a state-of-the-art randomized search heuristics in continuous domain [8–10, 12]. The CMA-ES maintains the Gaussian distribution, from which candidate solutions are drawn. It repeats the following: sample λ points from the Gaussian distribution, evaluate the fitness for each sample, update the parameters including the mean vector and the covariance matrix of the distribution in order to make the distribution likely to generate better solutions. Recently, it has been shown [2, 7] that the parameter update in the CMA-ES is partially interpreted as a natural gradient ascent on the parameter space of the Gaussian distribution, where the natural gradient is computed for the function defined below in (1). This idea is further generalized to the generic framework for arbitrary optimization, namely information-geometric optimization (IGO) [16].

Since the CMA-ES has been originally proposed for unconstrained continuous optimization, it often requires a treatment when solving a constrained problem. A number of constraint handling strategies have been proposed for evolution strategies and for more generic evolutionary algorithms [15]; e.g., adding an adaptive penalty to the fitness according to the constraint violation [11], repairing an infeasible point into the feasible region by a projection onto the boundary [4] or by a gradient based repair operator [13]. In this paper we consider the *resampling* strategy; an infeasible point is

[©] Springer International Publishing Switzerland 2014

discarded and resampled until it drops into the feasible region. It can be applied even when the constraint functions are black-box.

It has been shown that under a linear constraint, the success probability, i.e. the probability of generating a better point, depends on the angle of the gradients of the constraint function and the objective function and the dependency of the success probability on the angle is different for the resampling [5] and a repair operator [4]. This obviously affects a success probability based parameter update such as step-size adaptation based on the 1/5 success rule [17]. Moreover, since the distribution of the generated point in the feasible region is truncated (in the case of resampling) or biased on the boundary (in the case of repairing), the update rules that are designed from a statistical viewpoint are affected. For example, when the original Gaussian distribution is parameterised by the mean vector m and the covariance matrix C, these parameters no more represent the mean vector and the covariance matrix of the truncated Gaussian distribution. Then the maximum likelihood estimators for m and C for the truncated distribution differ from the ones for the original distribution. Therefore, a treatment in parameter update is required, an example of which is proposed by [6] where the covariance matrix is actively reduced in the direction of the gradient of the constraint.

In this paper we study the effect of the constraint from a viewpoint of the natural gradient. When resampling method is employed, the distribution of the generated feasible points is a truncated probability distribution whose domain is limited to the feasible set. In this situation the natural gradient differs from the one computed for the non-truncated (original) probability distribution. Now a question arises as to if we can gain a better performance by computing the natural gradient on the manifold of the truncated Gaussian distributions limited to the feasible set.

To address the question we derive the natural gradient under a linearly constrained feasible domain and compare it with the original natural gradient theoretically and numerically. In Section 2 the IGO framework and the rank- μ update CMA-ES as an instantiation of the IGO are revisited. In Section 3 we derive the natural gradient under a linearly constrained feasible domain. In Section 4 we analyze the infinite-population model using the derived natural gradient on a linearly constrained spherical problem and perform simulations to compare the behavior of the derived algorithm with the original algorithm on a linearly constrained spherical problem. Finally in Section 5 we summarize this work and discuss required future works.

Notation. The inner product of $x \in \mathbb{R}^d$ and $y \in \mathbb{R}^d$ is denoted by $\langle x, y \rangle$ and the norm of x by $||x|| = \langle x, x \rangle^{1/2}$. For any matrix A, $[A]_{i,j}$ represents the (i, j)th element, $[A]_{i,:}$ (or $[A]_{:,j}$) the *i*th row (or the *j*th column, respectively.) For any symmetric matrix A of dimension d, let vech(A) denote the lower-left half vectorization of A such that vech(A) is the d(d + 1)/2 dimensional column vector whose *i*th element is $[A]_{m_i,n_i}$ where $i = m_i + (d - n_i/2)(n_i - 1)$ for $1 \le n_i \le m_i \le d$. We refer to [14] for the detail.

2 IGO Framework and the Rank-μ Update CMA

Formulation. We consider a constrained continuous minimization $\operatorname{argmin}_{x \in X} f(x)$, where $X \subset \mathbb{R}^d$ is the feasible set and *f* is the objective function defined over X. In the

following sections we assume that the feasible set is restricted by a linear function, namely $X = \{x \in \mathbb{R}^d \mid \langle x, v \rangle \ge \alpha\}$ for some unit vector $v \in \mathbb{R}^d \setminus \{0\}$ and some $\alpha \in \mathbb{R}$.

Given a family \mathcal{P} of probability distributions P_{θ} on X parameterized by a real vector $\theta \in \Theta$, the IGO framework formulates the joint problem on the parameter space Θ at each iteration *t* as follows

$$\theta^{t+1} = \operatorname{argmax}_{\theta \in \Theta} J_{\theta^t}(\theta), \quad \text{where} \quad J_{\theta^t}(\theta) := \int_{\mathbf{X}} W^f_{\theta^t}(x) P_{\theta}(\mathrm{d}x) \quad .$$
(1)

Here θ^t is the value of the parameter at iteration *t*, $W_{\theta^t}^f$ defines the preference that is monotonic to *f*. The preference is defined based on the probability of sampling a better point; namely,

$$W_{\theta^{t}}^{f}(x) = w\left(P_{\theta^{t}}[y \in X \mid f(y) \leq f(x)]\right) \quad \text{, where } w : [0, 1] \to \mathbb{R}.$$
(2)

Another weight scheme is introduced in [1],

$$W_{\mu t}^{f}(x) = -\left(\mu_{\text{Leb}}[y \in X : f(y) \le f(x)]\right)^{2/d},$$
(3)

where μ_{Leb} denotes the Lebesgue measure on \mathbb{R}^d . This is theoretically attractive; on an unconstrained monotonic convex quadratic composite function $g(x^TAx)$ with g strictly increasing, this weight value is $-cx^TAx$, where c is a constant independent of g, m, and C, and it enables us to derive the exact $J_{\theta'}(\theta)$.

Natural Gradient. The *natural gradient* can be interpreted as the gradient of a function defined on the space of the probability distribution equipped with the Fisher metric. It can be also interpreted as the steepest ascent direction of the function with respect to the KL-divergence. Since the Fisher metric (and KL-divergence) is independent of the parameterization (coordinate system) of the probability distribution, the natural gradient is invariant to any re-parameterization of θ . Given a parameterization θ , the natural gradient is computed by the product of the inverse of the Fisher information matrix of θ and the vanilla gradient (gradient on the Euclidean space) of the log-likelihood of the probability distribution. We refer to [16] for further properties of the natural gradient.

Noting that $W^f_{\theta'}$ in (2) or (3) is independent of θ , the natural gradient of $J_{\theta'}$ is computed by

$$\tilde{\nabla} J_{\theta^{l}}(\theta) = \int_{Y} W^{f}_{\theta^{l}}(x) \tilde{\nabla} l(\theta; x) P_{\theta}(\mathrm{d}x) \quad , \tag{4}$$

where $\tilde{\nabla}$ represents the map from a function to its natural gradient, and $l(\theta; x) = \ln p_{\theta}(x)$ denotes the log-likelihood at θ given x. Eq (4) is viewed as a weighted expectation of the natural gradient of the log-likelihood at θ .

Implementation of the Natural Gradient Ascent. The IGO algorithm performs the natural gradient ascent instead of exactly solving joint problem (1). Then iterate $\{\theta^t\}$ is defined by

$$\theta^{t+1} = \theta^t + \eta^t \tilde{\nabla} J_{\theta^t}(\theta)|_{\theta = \theta^t} \quad , \tag{5}$$

where η^t denotes the step-size for the natural gradient ascent, aka the learning rate for the parameter update, which is sometimes replaced with a diagonal matrix whose diagonal entries are the learning rates for each element of the parameter vector. However, the integration in (4) cannot be performed analytically in advance unless *f* is known.

Therefore, we estimate (4) with samples x_1, \ldots, x_{λ} drawn from $P_{\theta'}$. According to [16], we can approximate $W^f_{\theta'}$ in (2) for each x_i as $W^f_{\theta'}(x_i) \approx \hat{w}_{\mathrm{rk}(x_i)} = w((\mathrm{rk}(x_i) - 1/2)/\lambda)$, where $\mathrm{rk}(x_i)$ denotes the ranking of $f(x_i)$ among $f(x_1), \ldots, f(x_{\lambda})$. With this, a Monte-Carlo estimate provides an approximation of (4) at $\theta = \theta^t$, namely

$$\tilde{\nabla} J_{\theta^{t}}(\theta)|_{\theta=\theta^{t}} \approx \frac{1}{d} \sum_{i=1}^{d} \hat{w}_{\mathrm{rk}(x_{i})} \tilde{\nabla} l(\theta^{t}; x_{i})$$
 (6)

The IGO implementation performs the natural gradient ascent (5) with replacing the natural gradient $\tilde{\nabla} J_{\theta'}(\theta)|_{\theta=\theta'}$ given in (4) with its approximation (6).

Rank-\mu Update CMA. Considering the IGO implementation for unconstrained continuous optimization, i.e. $X = \mathbb{R}^d$, with the Gaussian distributions on \mathbb{R}^d , it is known from [3] that the natural gradient of the log-likelihood of the Gaussian distribution is explicitly written in a special form. If the Gaussian distribution is parameterized by $\theta = [m^T, \operatorname{vech}(C)^T]$, where *m* and *C* are the mean vector and the covariance matrix, the parameter update in the IGO implementation reads

$$m^{t+1} = m^{t} + \frac{\eta_{m}}{\lambda} \sum_{i=1}^{\lambda} \hat{w}_{\mathrm{rk}(x_{i})}(x_{i} - m) C^{t+1} = C^{t} + \frac{\eta_{C}}{\lambda} \sum_{i=1}^{\lambda} \hat{w}_{\mathrm{rk}(x_{i})}((x_{i} - m)(x_{i} - m)^{\mathrm{T}} - C) .$$
(7)

This is called the rank- μ update [10] and is a component of the standard CMA [9].

In [3], $X = \mathbb{R}^d$ is assumed to obtain the explicit form for the natural gradient. In other words, the natural gradient computed in the reference is the one on the manifold of (nontruncated) Gaussian distributions defined on \mathbb{R}^d . If X is a proper subset of \mathbb{R}^d ($X \subset \mathbb{R}^d$ and $X \neq \mathbb{R}^d$) and the truncated Gaussian distribution is considered (sampling from a Gaussian distribution with resampling scheme as a constraint handling), the natural gradient on the manifold of the truncated Gaussian distributions on X is different from the one derived in [3] and the resulting natural gradient ascent differs from (7). This is the main concern of the paper.

3 Natural Gradient for Truncated Gaussian Distributions

Let $p_{\theta}(x)$ and $l(\theta; x)$ be the probability density function (p.d.f.) and the log-likelihood function (l.l.f.) induced by the Gaussian distribution P_{θ} with mean $m = m(\theta)$ and covariance matrix $C = C(\theta)$, i.e., $l(\theta; x) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln \det(C(\theta)) - \frac{1}{2}(x - m(\theta))^{\mathrm{T}}C^{-1}(\theta)(x - m(\theta))$ and $p_{\theta} = \exp(l(\theta; x))$. Then, $P_{\theta}(A) = \int_{A} p_{\theta}(x) dx$ for any Lebesgue measurable $A \subset \mathbb{R}^{d}$. As in the rank- μ update CMA-ES, we consider $\theta = [m^{\mathrm{T}}, \operatorname{vech}(C)^{\mathrm{T}}]^{\mathrm{T}}$.

If the X is a proper subset of \mathbb{R}^d and the resampling strategy is employed, the distribution of the samples in X is the Gaussian distribution truncated on X. Let $\bar{p}_{\theta}(x)$ and $\bar{l}(\theta; x)$ be the p.d.f. and l.l.f. of such a truncated Gaussian distribution \bar{P}_{θ} over X. Then, $\bar{p}_{\theta}(x) = p_{\theta}(x)/P_{\theta}(X)$ and $\bar{l}(\theta; x) = l(\theta; x) - \ln P_{\theta}(X)$ for $x \in X$, where $P_{\theta}(X) = \int_X p_{\theta}(x) dx = \mathbb{E}_{x \sim p_{\theta}}[\mathbb{I}\{x \in X\}]$ is the probability of x being sampled in X from p_{θ} . If we implement an algorithm following the IGO framework, the natural gradient on the manifold of the truncated Gaussian distributions $\{\bar{P}_{\theta} \mid \theta \in \Theta\}$ is needed.

As the first attempt of the work, we derive the natural gradient of the l.l.f. of the Gaussian distribution truncated on $X = \{x \in \mathbb{R}^d \mid \langle x, v \rangle \ge \alpha\}$. The following proposition and theorem provide the formula to compute the natural gradient explicitly.

Proposition 1. Let φ and Φ be the p.d.f. and cumulative density function induced by the standard normal distribution $\mathcal{N}(0, 1)$. Define φ_{β} be the p.d.f. for the normal distribution truncated onto $\{x \ge \beta\}$, that is, $\varphi_{\beta}(x) = \varphi(x)/(1 - \Phi(\beta))$. Let N_{β} be a random variable obeying φ_{β} . Then, $\mu_1 := \mathbb{E}[N_{\beta}] = \varphi_{\beta}(\beta), \mu_2 := \mathbb{E}[N_{\beta}^2] = \beta\varphi_{\beta}(\beta) + 1, \mu_3 := \mathbb{E}[N_{\beta}^3] = (\beta^2 + 2)\varphi_{\beta}(\beta)$ and $\mu_4 := \mathbb{E}[N_{\beta}^4] = (\beta^3 + 3\beta)\varphi_{\beta}(\beta) + 3$.

Theorem 1. Let the natural gradient $\tilde{\nabla}\overline{l}(\theta; x)$ of the l.l.f. of the truncated Gaussian distribution decomposed as $\tilde{\nabla}\overline{l}(\theta; x) = [\delta n(x)^{\mathrm{T}}, \operatorname{vech}(\delta C(x))^{\mathrm{T}}]^{\mathrm{T}}$, where $\delta n(x) \in \mathbb{R}^d$ and $\delta C(x) \in \mathbb{R}^{d \times d}$ are the components corresponding to m and C respectively. Let $u = v/|C^{1/2}v||$ and define y = x - m. Then,

$$\begin{split} \delta m(x) &= \left[\frac{\mu_2 - \mu_1 \langle u, y \rangle}{\tau_1}\right] y + \left[-\left(\frac{\tau_2}{\tau_1 \tau_3 - \tau_2^2} - \frac{\mu_1}{\tau_1}\right) \langle u, y \rangle^2 \\ &+ \left(\frac{\tau_3}{\tau_1 \tau_3 - \tau_2^2} - \frac{\mu_2}{\tau_1}\right) \langle u, y \rangle - \frac{\tau_3 \mu_1 - \tau_2 \mu_2}{\tau_1 \tau_3 - \tau_2^2} \right] C u \quad and \end{split} \tag{8}$$
$$\delta \mathcal{C}(x) &= y y^{\mathrm{T}} - C + \left[\frac{(1 - \tau_1) \langle u, y \rangle + \mu_1}{\tau_1}\right] (y u^{\mathrm{T}} C + C u y^{\mathrm{T}}) + \left[\left(2 \frac{\tau_2 \mu_1 - \tau_1 \mu_2}{\tau_1 \tau_3 - \tau_2^2} + 1\right) \right. \\ &+ \left(\frac{2\tau_1}{\tau_1 \tau_3 - \tau_2^2} - \frac{2}{\tau_1} + 1\right) \langle u, y \rangle^2 - 2 \left(\frac{\tau_2}{\tau_1 \tau_3 - \tau_2^2} + \frac{\mu_1}{\tau_1}\right) \langle u, y \rangle \right] C u u^{\mathrm{T}} C \ , \tag{9}$$

where μ_1, μ_2, μ_3 and μ_4 are as defined in Proposition 1 with $\beta = (\alpha - \langle v, m \rangle)/||C^{1/2}v||$, and $\tau_1 = \mu_2 - \mu_1^2, \tau_2 = \mu_3 - \mu_1\mu_2, \tau_3 = \mu_4 - \mu_2^2$.

We have omitted its proof due to the space limitation. Comparing to the natural gradient $\tilde{\nabla}l(\theta; x)$ of the l.l.f. for the non-truncated Gaussian distribution P_{θ} that can be expressed as $\delta n(x) = y$ and $\delta C(x) = yy^{T} - C$, (8) and (9) have additional components characterized by $Cu = Cv/||C^{1/2}v||$. The coefficients are determined by β —a signed distance to the boundary normalized by the standard deviation $||C^{1/2}v||$ in the direction of v—and $\langle u, y \rangle$ —a signed distance from the current mean m to the sample point x in the direction of $C^{1/2}v$. In the limit $\beta \to -\infty$, meaning that the constraint boundary is far away from the current mean and the situation is close to the unconstrained case, we have from Proposition 1 that $\mu_1 = \mu_3 = \tau_2 = 0$, $\mu_2 = \tau_1 = 1$, $\mu_4 = 3$, $\tau_3 = 2$, and (8) and (9) become identical to the natural gradient for the unconstrained case.

The natural gradient of the l.l.f. only depends on the manifold of the probability distributions. That is, it only depends on the feasible set X, but not on the objective function f. The parameter update (5) with (6) on the other hand depends on the selection scheme. More precisely, the adjustments δn and δC of the parameters is the weighted sum of $\delta n(x_i)$ and $\delta C(x_i)$ over $i = 1, ..., \lambda$, where the weight value is determined by the ranking of $f(x_i)$. In the next section we demonstrate on a linearly constrained spherical problem how much the derived natural gradient differs from the rank- μ update (7).

4 Study on a Linearly Constrained Spherical Problem

We consider the following linearly constrained spherical problem $\operatorname{argmin}_{x \in X_{\alpha}} f(x) := g(||x||^2)$, where $X_{\alpha} = \{x \in \mathbb{R}^d \mid \langle x, v \rangle \ge \alpha\}$ and g is strictly increasing. If $\alpha \ge 0$, the optimum is located on the boundary $x^* = \alpha v$, otherwise $x^* = (0, \dots, 0)$ and the landscape

around the optimum is the same as the unconstrained sphere function. Therefore, we consider only $\alpha \ge 0$ in this work.

For ranking-based weight scheme as in (2), the natural gradient (4) generally needs to be approximated by (6). To understand and emphasize the benefit of use the natural gradient derived in the previous section, we employ the weight scheme (3). As mentioned after (3), we can compute the joint objective analytically, $J_{\theta'}(\theta) = -c(||m||^2 + \text{Tr}(C))$, on the unconstrained spherical problem and the natural gradients become $\delta n = -2cCm$ and $\delta C = -2cC^2$ with an appropriate constant factor *c*. This analytical natural gradient is the limit of the natural gradient estimate (6) w.r.t. $\lambda \to \infty$ [1] and it models the infinite-population behavior of the rank- μ update CMA (7).

If $\alpha = 0$, the volume of each sub level set $\mu_{\text{Leb}}[y \in X : f(y) \leq f(x)]$ is just halved compared to the unconstrained case and we still have similar results.

Lemma 1. If $\alpha = 0$, the weight $W_{\theta^{f}}^{f}(x)$ defined in (3) is $-\tilde{c}x^{T}x$, where \tilde{c} is a constant independent of m and C. The natural gradient (4) with δm and δC derived in Theorem 1 reads $\delta m = -2\tilde{c}Cm$ for m and $\delta \overline{C} = -2\tilde{c}C^{2}$ for C.

Surprisingly, the natural gradient on the linearly constrained spherical problem is only different in length from the one on unconstrained spherical problem. This implies that the natural gradient update (5) with δn and δC in Theorem 1 reads the exact same parameter update as in the unconstrained case with an appropriate η^t . Therefore, all the results in [1] are carried over here as stated in the next theorem.

Theorem 2. Let λ_1^t denote the largest eigenvalue of $-C^{-1/2}\overline{\partial C}C^{-1/2}$. If C^0 is symmetric positive definite and $\eta^t \lambda_1^t < 1$ for all $t \ge 0$, then C^t is symmetric positive definite. Moreover, if $\eta^t = \overline{\eta}/\lambda_1^t$ for $\overline{\eta} \in (0, 1/2]$, $\lim_{t\to\infty} \text{Cond}(C^t) = 1$ and $\lim_{t\to\infty} ||C^{t+1}||_F/||C^t||_F = \lim_{t\to\infty} ||m^{t+1}||/||m^t|| = 1 - \overline{\eta}$, where $||\cdot||_F$ denotes the Frobenius norm.

The learning rate $\eta^t = \bar{\eta}/\lambda_1^t$ is taken from [1]. This theorem means, the *C* becomes proportional to the Hessian of $x^T x$, namely the identity matrix, and *m* linearly converges towards the global optimum at the origin. For the detail, see [1].

To visualize the difference from the original parameter update (7) with $\lambda = \infty$ where the adjustment is $\mathbb{E}[W^f_{\theta^f}(x)(x-m)]$ and $\mathbb{E}[W^f_{\theta^f}(x)((x-m)(x-m)^T - C)]$, we derive the explicit form of the expectation. Following proposition reads it when the weight scheme with baseline subtraction, $W^f_{\theta^f}(x) - \mathbb{E}[W^f_{\theta^f}(x)]$, is introduced.

Proposition 2. Let y = x - m, $u = v/||C^{1/2}v||$ and μ_1 , μ_2 , τ_1 , τ_2 and τ_3 be as appeared in Theorem 1. If $\alpha = 0$, $\mathbb{E}[W^f_{\theta}(x)] = -\tilde{c}[\operatorname{Tr}(C) + ||m||^2 + 2\mu_1 u^{\mathrm{T}} Cm + (\mu_2 - 1)u^{\mathrm{T}} C^2 u]$, $\mathbb{E}[y] = \mu_1 Cu$, $\mathbb{E}[yy^{\mathrm{T}} - C] = (\mu_2 - 1)Cuu^{\mathrm{T}}C$, and

$$\mathbb{E}[(W_{\theta^{f}}^{f}(x) - \mathbb{E}[W_{\theta^{f}}^{f}(x)])y] = -\tilde{c}[(\tau_{2} - 2\mu_{1})(u^{T}C^{2}u) + 2(\tau_{1} - 1)(u^{T}Cm)]Cu - 2\tilde{c}\mu_{1}C^{2}u - 2\tilde{c}Cm \quad and \quad (10)$$

$$\mathbb{E}[(W_{\theta^{f}}^{f}(x) - \mathbb{E}[W_{\theta^{f}}^{f}(x)])(yy^{T} - C)] = -2\tilde{c}C^{2} - \tilde{c}[(\tau_{3} - 4\mu_{2} + 2)(u^{T}C^{2}u) + 2(\tau_{2} - 2\mu_{1})(u^{T}Cm)]Cuu^{T}C - 2\tilde{c}(\mu_{2} - 1)(C^{2}uu^{T}C + Cuu^{T}C^{2}) - 2\tilde{c}u_{1}(Cmu^{T}C + Cum^{T}C) . \quad (11)$$

Note that the $\mathbb{E}[W_{\theta'}^f(x)y] = \mathbb{E}[(W_{\theta'}^f(x) - \mathbb{E}[W_{\theta'}^f(x)]y] + \mathbb{E}[W_{\theta'}^f(x)]\mathbb{E}[y]$ and $\mathbb{E}[W_{\theta'}^f(x)(yy^{\mathrm{T}} - C)] = \mathbb{E}[(W_{\theta'}^f(x) - \mathbb{E}[W_{\theta'}^f(x)])(yy^{\mathrm{T}} - C)] + \mathbb{E}[W_{\theta'}^f(x)]\mathbb{E}[yy^{\mathrm{T}} - C]$. We call them NG_n



Fig. 1. Transitions of $||m||^2$, m_1 , the eigenvalues of *C*, the condition number Cond(*C*), and β from left to right, respectively. For NG_n, the run stops after β reaches 7. In the center figure, the eigenvalue corresponds to the first coordinate is the larger one for NG_b, and the smaller one for NG_n.

(natural gradient computed on the non-truncated Gaussian manifold), in contrast to NG_t (natural gradient computed on the truncated Gaussian manifold) for δn and δC derived in Lemma 1. Moreover, we call (10) and (11) NG_b. The only difference between NG_n and NG_b is the offset of the weight, $-\mathbb{E}[W_{\theta'}^f(x)]$, and the expectation of the weight in NG_b is forced to be zero. Note that this offset does not affect the natural gradient in NG_t derived Lemma 1 since the expectations of $\delta n(x)$ and $\delta C(x)$ taken over x are zero.

Fig. 1 illustrates the evolution of the parameters *m* and *C* following the natural gradient update (5) with the natural gradient $\tilde{\nabla} J_{\theta'}(\theta) = [\overline{\delta m}^{\mathrm{T}}, \operatorname{vech}(\overline{\delta C})^{\mathrm{T}}]^{\mathrm{T}}$ where $\overline{\delta m}$ and $\overline{\delta C}$ are computed for NG_t, NG_n, and NG_b. For the constraint we set $v = e_1 = [1, 0, \dots, 0]^{\mathrm{T}}$ and $\alpha = 0$. The step-size is $\eta^t = \overline{\eta}/\lambda_1^t$, where $\overline{\eta} = 0.1$ and λ_1^t is the largest eigenvalue of $-C^{-1/2}\overline{\delta C}C^{-1/2}$ with corresponding $\overline{\delta C}$ for each variant. This step-size setting guarantees the positivity of the covariance matrix as stated in Theorem 2 for NG_t. The problem dimension is d = 10. To produce simple figures, the evolution starts from $m^0 = e_1$ and $C^0 = I$. Thanks to the symmetry, *m* stays on the first axis and *C* stays to be a diagonal matrix whose second to the last diagonal elements are equal.

As stated in Theorem 2, *m* and *C* converge linearly in NG_t while the condition number of *C* stays 1 forever. In contrast, *m* goes over the constraint boundary and tends to stop at some point in the infeasible area while the condition number of *C* grows up in NG_b and NG_n. The normalized and signed distance β from *m* to the constraint boundary then becomes large, which in the actual algorithms means that the probability of sampling a point in the feasible region decreases. Since *m* does not converge towards the optimum, the best-so-far point would not converge linearly towards the global optimum. The tendency of the plots does not depend on the choice of the learning rate, i.e. β does not converge to zero with any learning rate in NG_n and NG_b.

So far the natural gradient is analytically computed. This is considered the approximated behavior of the algorithm in the limit of $\lambda \to \infty$. In practice the population size $\lambda < \infty$ and $W^f_{\theta'}(x)$ for each x_i and then $\overline{\delta n}$ and $\overline{\delta C}$ must be estimated using finite samples $x_1, \ldots, x_\lambda \sim \overline{p}_{\theta'}(x)$.¹ We can approximate $W^f_{\theta'}(x)$ as $\hat{W}^f_{\theta'}(x) :=$

¹ The resampling can be performed efficiently as follows. Generate $z \sim \mathcal{N}(0, I_d)$. If $\langle z, C^{1/2}u \rangle < \beta$, generate $\tilde{z} \sim \mathcal{N}(0, 1)$, resample it till $\tilde{z} \ge \beta$, then update $z = z + (\tilde{z} - \langle z, C^{1/2}u \rangle)C^{1/2}u$. Then, $m + C^{1/2}z$ obeys \bar{p}_{θ} . So we only need to resample a standard normal random number \tilde{z} .



Fig. 2. Transitions of $||m||^2 - m_1^2$, m_1^2 , the eigenvalues of *C*, the condition number Cond(*C*), and β for NG_t and the rank- μ update CMA with $\eta_m = \eta_C = \bar{\eta} = 0.1$ (above) and $\eta_m = \eta_C = \bar{\eta} = 0.01$ (below). The graphs are the average over 30 trials.

 $\begin{bmatrix} \sum_{j \in \{k \in [1]; \lambda] \mid f(x_k) \leq f(x)\}} (1/\bar{p}_{\theta'}(x_j)) \end{bmatrix}^{2/d}, \text{ where } \bar{p}_{\theta'}(x) = p_{\theta'}(x)/P_{\theta'}(X) \text{ and } P_{\theta'}(X) = 1 - \Phi(\beta) = [1 - \operatorname{erf}(\beta/\sqrt{2})]/2.^2 \text{ With this approximation we can estimate the weight with baseline subtraction}^3, W_{\theta'}^f(x) - \mathbb{E}[W_{\theta'}^f(x)], \text{ for each } x_i \text{ as } w_i = \hat{W}_{\theta'}^f(x_i) - \frac{1}{\lambda} \sum_{j=1}^{\lambda} \hat{W}_{\theta'}^f(x_j).$ Then, $\overline{\delta n}$ and $\overline{\delta C}$ are approximated by the average of $w_i \cdot \delta n(x_i)$ and $w_i \cdot \delta C(x_i)$, respectively. We denote the estimated natural gradient by $\widehat{\delta n}$ and $\widehat{\delta C}$.

Fig. 2 compares the behaviors of NG_t with that of the rank- μ update CMA (7) (denoted CMA) with the weight used in the standard CMA, $\hat{w}_{\text{rk}(x_i)} = \max(0, \ln(\frac{\lambda+1}{2}) - \ln(rk(x_i))) / \sum_{j=1}^{\lambda} \max(0, \ln(\frac{\lambda+1}{2}) - \ln(j))$. The problem dimension d = 10, the population size $\lambda = 1000$, and the learning rates $\eta_m = \eta_C = 0.1$ and 0.01 for the rank- μ update CMA. For NG_t, $\eta^t = \bar{\eta} / \max[\sigma(C^{-1/2}\widehat{\delta C}C^{-1/2}), ||C^{-1/2}\widehat{\delta n}||]$ with $\bar{\eta} = 0.1$ and 0.01, where $\sigma(\cdot)$ represents the largest singular value. The graphs show the average of 30 independent runs for each method.

As you can see from the figure, in the rank- μ update CMA, β tends to stay at some point in negative, meaning that the mean vector is always away from the constraint boundary in the feasible domain and its distance is proportional to $||C^{1/2}v||$. The eigenvalue of *C* corresponding to $v = e_1$ becomes relatively smaller than the other eigenvalues.

² In practice, a scalar factor in $\hat{W}_{\theta^{t}}^{f}$ does not matter at all because the natural gradient is multiplied by η_{m}^{t} and η_{C}^{t} that are inversely proportional to the scalar factor as introduced below. Therefore, we can replace $1/\bar{p}_{\theta^{t}}(x)$ with $\exp(||z||^{2}/2)$, where $z = C^{-1/2}(x - m)$.

³ As stated above, the baseline subtraction in NG_t does not affect the expectation of the natural gradient, while it can reduce the estimation variance of the natural gradient.

On the other hand, NG_t results in smaller values of β and Cond(*C*) and they get even smaller if we decrease the learning rate or increase the population size.

5 Summary and Discussion

In this paper we derive the natural gradient of the l.l.f. of the truncated Gaussian distribution for linearly constrained optimization problems. Analysis on a linearly constrained spherical problem shows the infinite-population model using the derived natural gradient reads the same update as the exact natural gradient algorithm on a unconstrained spherical problem [1] and all the results proven in the reference hold. The simulation results exhibit different behavior of the derived algorithm and the rank- μ update CMA. The rank- μ update CMA tends to stay in the feasible set and the distance from the constraint boundary stays proportional to $||C^{1/2}v||$, where v is the normal vector of the constraint boundary, whereas the condition number and the normalized distance in the derived algorithm converges to smaller values than in the rank- μ update CMA.

We would like to remark that the simulation performed in Section 4 depends heavily on the weight scheme. As we see in Fig. 1 and Fig. 2, NG_n, NG_b, and the rank- μ update CMA result in different behavior, although their only difference is the weight value and the learning rate. Moreover, from a preliminary experiment we have observed that NG_t does not work as well as it is with the weight scheme (3) if we employ the CMA-type weight scheme or the fitness proportional weight $\hat{W}_{\theta}^{f}(x_{i}) = ||x||^{2}$. The mean vector enters the infeasible region as we have observed in NG_n. Especially for the fitness proportional weight, we can derive a theoretical result for the infinite-population model that even if $\alpha > 0$, the natural gradient becomes the same as the one on the unconstrained sphere problem and the mean vector tends to converges to the origin that is in the infeasible domain. Further study on the weight scheme is highly required.

For other future works, we compare the derived algorithm with the existing treatment for the constrained problem such as [6]. To enhance the performance, we would need to incorporate a step-size control mechanism that is in general heavily affected by the constraint, and a projection of the mean vector to the feasible domain when it reaches the infeasible domain. Furthermore, we extend the formula for the natural gradient for a linearly constrained problem stated in Theorem 1 to problems with more general constraint.

Acknowledgments. This work is supported by JSPS KAKENHI Grant Number 25880012.

References

- Akimoto, Y.: Analysis of a Natural Gradient Algorithm on Monotonic Convex-Quadratic-Composite Functions. In: Genetic and Evolutionary Computation Conference, pp. 1293–1300 (2012)
- Akimoto, Y., Nagata, Y., Ono, I., Kobayashi, S.: Bidirectional relation between CMA evolution strategies and natural evolution strategies. In: Schaefer, R., Cotta, C., Kołodziej, J., Rudolph, G. (eds.) PPSN XI. LNCS, vol. 6238, pp. 154–163. Springer, Heidelberg (2010)

- Akimoto, Y., Nagata, Y., Ono, I., Kobayashi, S.: Theoretical Foundation for CMA-ES from Information Geometry Perspective. Algorithmica 64, 698–716 (2012)
- Arnold, D.V.: Analysis of a repair mechanism for the (1, λ)-ES applied to a simple constrained problem. In: Genetic and Evolutionary Computation Conference, pp. 853–860 (2011)
- 5. Arnold, D.V.: On the behaviour of the $(1, \lambda)$ -ES for a simple constrained problem. In: Foundations of Genetic Algorithms, pp. 15–24 (2011)
- Arnold, D.V., Hansen, N.: A (1 + 1)-CMA-ES for constrained optimisation. In: Genetic and Evolutionary Computation Conference Conference, pp. 297–304 (2012)
- Glasmachers, T., Schaul, T., Sun, Y., Wierstra, D., Schmidhuber, J.: Exponential Natural Evolution Strategies. In: Genetic and Evolutionary Computation Conference, pp. 393–400 (2010)
- Hansen, N.: Benchmarking a BI-population CMA-ES on the BBOB-2009 noisy testbed. In: Companion on Genetic and Evolutionary Computation Conference (2009)
- Hansen, N., Auger, A.: Principled Design of Continuous Stochastic Search: From Theory to Practice. In: Borenstein, Y., Moraglio, A. (eds.) Theory and Principled Methods for the Design of Metaheuristics. Springer (2013)
- Hansen, N., Muller, S.D., Koumoutsakos, P.: Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES). Evolutionary Computation 11(1), 1–18 (2003)
- Hansen, N., Niederberger, A.S.P., Guzzella, L., Koumoutsakos, P.: A Method for Handling Uncertainty in Evolutionary Optimization With an Application to Feedback Control of Combustion. IEEE Transactions on Evolutionary Computation 13(1), 180–197 (2009)
- Hansen, N., Ostermeier, A.: Completely derandomized self-adaptation in evolution strategies. Evolutionary Computation 9(2), 159–195 (2001)
- Harada, K., Sakuma, J., Ono, I., Kobayashi, S.: Constraint-handling method for multiobjective function optimization: Pareto descent repair operator. In: Obayashi, S., Deb, K., Poloni, C., Hiroyasu, T., Murata, T. (eds.) EMO 2007. LNCS, vol. 4403, pp. 156–170. Springer, Heidelberg (2007)
- 14. Harville, D.A.: Matrix Algebra from a Statistician's Perspective. Springer (2008)
- 15. Kramer, O.: A review of constraint-handling techniques for evolution strategies. Applied Computational Intelligence and Soft Computing 2010 (2010)
- Ollivier, Y., Arnold, L., Auger, A., Hansen, N.: Information-geometric optimization algorithms: A unifying picture via invariance principles (2011), http://arxiv.org/abs/1106.3708
- Rechenberg, I.: Evolutionsstrategie: Optimierung technischer Systeme nach Prinzipien der biologischen Evolution. Frommann-Holzboog (1973)