

## SMART GRADIENT - AN ADAPTIVE TECHNIQUE FOR IMPROVING GRADIENT ESTIMATION

ESMAIL ABDUL FATTAH\*, JANET VAN NIEKERK AND HÅVARD RUE

CEMSE Division, King Abdullah University of Science and Technology, Kingdom of Saudi Arabia

**ABSTRACT.** Computing the gradient of a function provides fundamental information about its behavior. This information is essential for several applications and algorithms across various fields. One common application that requires gradients are optimization techniques such as stochastic gradient descent, Newton's method and trust region methods. However, these methods usually require a numerical computation of the gradient at every iteration of the method which is prone to numerical errors. We propose a simple limited-memory technique for improving the accuracy of a numerically computed gradient in this gradient-based optimization framework by exploiting (1) a coordinate transformation of the gradient and (2) the history of previously taken descent directions. The method is verified empirically by extensive experimentation on both test functions and on real data applications. The proposed method is implemented in the **R** package `smartGrad` and in **C++**.

**1. Introduction.** Gradients are one of the oldest constructs of modern mathematics and often form the basis for many optimization problems. Moreover, gradients are used in various function expansions like the Taylor series expansion, in tangent plane construction and in various real-life engineering challenges such as building ramps, roads and buildings, amongst others. The gradient of a mathematical function  $f$ , at the input  $\mathbf{x}$ , is denoted by  $\nabla f(\mathbf{x})$ . It can be mathematically interpreted as a rate of disposition based on the disposition of  $\mathbf{x}$  or graphically as the slope of the tangent plane at  $\mathbf{x}$  [6]. Often though, the analytical form of  $\nabla f(\mathbf{x})$  is unknown or computationally intensive to evaluate and hence the popularity of numerical gradients methods. Newton's and quasi-Newton methods [3] are well-known frameworks in optimization that use (numerical) gradients to get the descent directions.

Consider the problem of minimizing a twice differentiable continuous function  $f$ , where  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ , is convex, i.e. satisfies,

$$f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y})$$

for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and  $\lambda \in [0, 1]$ . Assume this (unconstrained) optimization problem is solvable, where its optimum  $\mathbf{x}^*$  exists and is unique. In this scenario, a necessary and a sufficient condition for  $\mathbf{x}^*$  to be optimal is

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

---

2020 *Mathematics Subject Classification.* Primary: 03-08; Secondary: 62-08.

*Key words and phrases.* Adaptive Technique, Gradient Estimation, Numerical Gradient, Optimization, Vanilla Gradient Descent.

\* Corresponding author: Esmail Abdul Fattah, email: [esmail.abdulfattah@kaust.edu.sa](mailto:esmail.abdulfattah@kaust.edu.sa).

Here,  $\nabla f(\mathbf{x})$  is composed of the function's partial derivatives, that describe the rates of change in multiple dimensions. More generally we consider the directional derivative of a function. Given a continuous function  $g$  with its first order partial derivatives, the directional derivative [10] of  $g$  at  $\mathbf{x} \in \mathbb{R}^n$  in the direction of unit vector  $\mathbf{u}$  is

$$D_{\mathbf{u}}g(\mathbf{x}) = \frac{\partial g}{\partial x_1}u_1 + \frac{\partial g}{\partial x_2}u_2 + \dots + \frac{\partial g}{\partial x_n}u_n.$$

The gradient vector can be reformulated as a combination of directional derivatives based on the canonical basis  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ , where every element of  $\mathbf{e}_i$  is zero except for the  $i^{\text{th}}$  position being 1, as in formula (1), that will follow shortly.

Common methods for solving optimizations problems rely on iterative techniques which generate iterations that converge to the optimum  $\mathbf{x}^*$ . The iterations are constructed from the general iterative update equation,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

with step size  $\alpha_k$  in the direction of the vector  $\mathbf{d}^{(k)}$ . The choice of the directions  $\{\mathbf{d}_k\}_{k=1}^n$  depends on the utilized method [5]:

1. Gradient Descent:  $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$ .
2. Quasi Newton's Method:  $\mathbf{d}^{(k)} = -\mathbf{B}_k \nabla f(\mathbf{x}^{(k)})$  where  $\mathbf{B}_k$  is an estimate of the Hessian at  $\mathbf{x}^{(k)}$ ,  $\nabla^2 f(\mathbf{x}^{(k)})$ .
3. Newton's Method:  $\mathbf{d}^{(k)} = -\nabla^2 f(\mathbf{x}^{(k)})^{-1} \nabla f(\mathbf{x}^{(k)})$ .

The choice of the step size  $\{\alpha^{(k)}\}$  is obtained using either exact or inexact line search [6]. Although second order methods are already good enough as they utilize Hessian information, here we focus on improving the first order methods such as gradient descent and quasi-newton methods that include symmetric rank 1 method and BFGS [3].

We summarize and relate four popular gradient computational frameworks next [2].

1. **Exact gradient with the canonical basis:** Based on directional derivatives, the gradient can be computed using the canonical basis  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ ,

$$\nabla f(\mathbf{x}) = \left( D_{\mathbf{e}_1} f(\mathbf{x}), D_{\mathbf{e}_2} f(\mathbf{x}), \dots, D_{\mathbf{e}_n} f(\mathbf{x}) \right). \quad (1)$$

2. **Exact gradient with a non-canonical basis:** Here, the gradient is computed based on the directional derivatives using a set of directions obtained from  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ . For a non-singular matrix  $\mathbf{G} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_n]$ , then

$$\nabla_{\mathbf{v}} f(\mathbf{x}) = \mathbf{G}^{-T} \nabla h(\boldsymbol{\varphi}) \Big|_{\boldsymbol{\varphi}=0} \text{ where } h(\boldsymbol{\varphi}) = f(\mathbf{x} + \mathbf{G}\boldsymbol{\varphi}). \quad (2)$$

3. **Inexact gradient in canonical basis: Vanilla Gradient (VG).** The gradient of the objective function is computed numerically as an estimate to the exact gradient with the canonical basis, such as using finite-difference methods,

$$\tilde{\nabla} f(\mathbf{x}) \approx \nabla f(\mathbf{x}). \quad (3)$$

$$\tilde{\nabla} f(\mathbf{x}) = \begin{cases} \text{forward finite difference,} \\ \text{backward finite difference,} \\ \text{central finite difference,} \\ \vdots \end{cases}$$

4. **Inexact gradient in non-canonical basis:** The gradient is computed using a general basis instead of the canonical basis,

$$\tilde{\nabla}_v f(\mathbf{x}) = \mathbf{G}^{-T} \tilde{\nabla} h(\boldsymbol{\varphi}) \Big|_{\boldsymbol{\varphi}=0} \approx \nabla f(\mathbf{x}) \text{ where } h(\boldsymbol{\varphi}) = f(\mathbf{x} + \mathbf{G}\boldsymbol{\varphi}). \quad (4)$$

For  $\mathbf{G}^T = \mathbf{I}_n$ , (2) reduces to (1), and (4) reduces to the Vanilla Gradient in (3) where the partial derivatives are directly computed along the vectors of the canonical basis.

A natural question is if we can construct  $\mathbf{G}$  in (4) that results in a more accurate Vanilla Gradient. As an illustration, we estimate the gradient of the two-dimensional Rosenbrock function [1]  $f(\mathbf{x}) = (1 - x_1)^2 + 100 \times (x_2 - x_1^2)^2$ , see Figure 1, using various bases. We use the central difference method of first order of step length  $10^{-3}$  to get the gradient estimate at  $\mathbf{x}^{(0)} = (-0.29, 0.4)^T$ . As for the directions, we start with the unit vectors  $\mathbf{e}_1 = (1, 0)^T$  and  $\mathbf{e}_2 = (0, 1)^T$  then we keep rotating with an angle of  $\pi/10^3$  until most of the directions are explored. The mean square error for the estimated gradient is calculated for each set of directions, and at each direction the magnitude of the gradient is calculated. The results are presented in Figure 1.

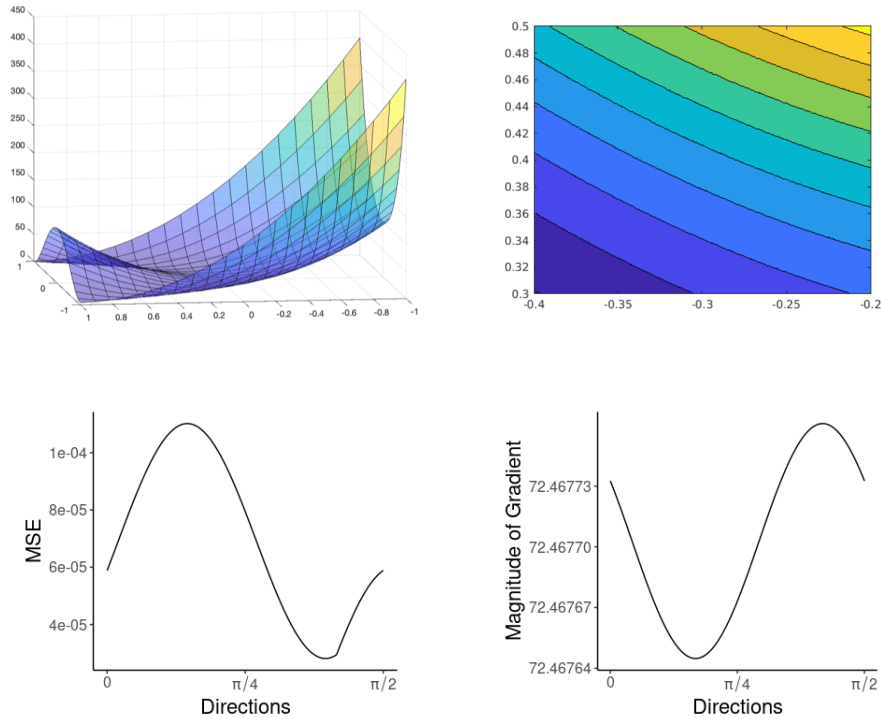


FIGURE 1. 2-dimensional Rosenbrock Function and a contour plot around the point  $p(-0.29, 0.4)$ . At the same point  $p$ , the MSE and magnitude of gradient are plotted using different directions.

From Figure 1, we observe opposite trends between the error in the gradient estimates and the magnitude of the gradient itself: the MSE in the gradient estimate is low when the magnitude of the gradient is high. Although such a claim doesn't hold in general, we found empirically it holds on most functions we tested for the

average MSE throughout iterations, and this is sufficient for our proposed method. This observation, suggests the obvious approach, here explained in dimension 2 for simplicity, to improve the gradient estimates at  $\mathbf{x}$ ,

1. Find the direction  $\mathbf{d}$  that maximize the change in the function itself.
2. Estimate the gradient in direction  $\mathbf{d}$  and the direction orthogonal to  $\mathbf{d}$ .
3. Do a linear transformation for the estimates we get in 2 to obtain the estimates of the Vanilla Gradient as in (4).

We are moving around the same point: to get the descent direction we estimate the gradient and to estimate the gradient we use direction where there is a high change in the objective value  $f$ , and this reduction is exactly what we want for the descent direction. This means that the directions we should use for estimating the gradient are exactly what we want to obtain by computing gradients. Within an iterative optimization algorithm, we have access to  $\mathbf{x}$  at most recent iterations,  $\mathbf{x}^{(k)}, \mathbf{x}^{(k-1)}, \dots$ , then the most recent change in  $\mathbf{x}$  can be used as the (surrogate) direction  $\mathbf{d}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$ . In dimension  $n$ , we would need to use the  $n$  most recent differences as surrogate directions. These directions should still be relevant, and hence we assume a low or moderate dimension of  $\mathbf{x}$ .

In this paper, we propose in Section 2 a method for constructing the matrix  $\mathbf{G}$  adaptively. In Section 3, we demonstrate how this construction results in improving the accuracy of a numerical gradient using test functions and Section 4 is for some applications including the R-INLA package. The paper is concluded by a discussion in Section 5.

**2. Methodology.** The proposed method is based on using the prior information we have about the descent directions in previous iterations,  $\{\mathbf{x}^{(i)}\}_{i=k-n}^k$  for improving the estimated numerical gradient  $\tilde{\nabla}f(\mathbf{x})$ . The difference between  $\mathbf{x}^{(i)}$  and  $\mathbf{x}^{(i-1)}$ ,  $i = k, \dots, k-n$ , indicates the direction where the objective function is highly reduced at that position. We then use these  $n$  differences as (surrogate) directions to estimate the gradient.

At iteration  $k$ , the  $n$  most recent differences can be used and the  $n$  directions have the following form,

$$\mathbf{d}^{(k)} = \frac{\Delta\mathbf{x}^{(k)}}{\|\Delta\mathbf{x}^{(k)}\|}$$

where  $\Delta\mathbf{x}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$ . To get a unique contribution from each direction  $\mathbf{d}^{(k)}$ , we subtract all the components of  $\mathbf{d}^{(k)}$  that are in the same direction of  $\{\mathbf{d}^{(j)}\}_{j=k-n+1}^{k-1}$ , then normalize it. We can do this by expressing these directions in terms of projectors,

$$\tilde{\mathbf{d}}^{(k)} = \mathbf{P}_{\perp\tilde{\mathbf{d}}^{(k-1)}}\mathbf{P}_{\perp\tilde{\mathbf{d}}^{(k-2)}}\dots\mathbf{P}_{\perp\tilde{\mathbf{d}}^{(k-n+1)}}\mathbf{d}^{(k)}$$

where  $\mathbf{P}_{\perp\tilde{\mathbf{d}}^{(k)}} = \mathbf{I} - \tilde{\mathbf{d}}^{(k)}(\tilde{\mathbf{d}}^{(k)})^T$ . We apply this orthogonalization to the  $n$  given directions using the Modified Gram-Schmidt (MGS) orthogonalization [7], which is an algorithm used to compute an orthonormal basis  $\{\tilde{\mathbf{d}}_1^{(k)}, \tilde{\mathbf{d}}_2^{(k)}, \dots, \tilde{\mathbf{d}}_n^{(k)}\}$  that spans the same subspace as the original vectors, i.e,

$$\text{Span}(\{\mathbf{d}_1^{(k)}, \mathbf{d}_2^{(k)}, \dots, \mathbf{d}_n^{(k)}\}) = \text{Span}(\{\tilde{\mathbf{d}}_1^{(k)}, \tilde{\mathbf{d}}_2^{(k)}, \dots, \tilde{\mathbf{d}}_n^{(k)}\}).$$

Some tiny noise to the  $\mathbf{x}$  differences was added to avoid singularities. This is as an adjustment when at least two directions are identical.

The matrix  $\tilde{\mathbf{G}}^{(k)}$  represents the orthogonal directions at iteration  $k$ ,

$$\tilde{\mathbf{G}}^{(k)} = [\tilde{\mathbf{d}}_1^{(k)} | \tilde{\mathbf{d}}_2^{(k)} | \dots | \tilde{\mathbf{d}}_n^{(k)}],$$

and is the MGS of matrix  $\mathbf{G}^{(k)}$ ,

$$\mathbf{G}^{(k)} = [\Delta \mathbf{x}^{(k)} | \mathbf{d}_1^{(k-1)} | \mathbf{d}_2^{(k-1)} | \dots | \mathbf{d}_{n-1}^{(k-1)}].$$

We use matrix  $\tilde{\mathbf{G}}^{(k)}$  to estimate the gradient at  $\mathbf{x}^{(k)}$ , using (4),

$$\tilde{\nabla}_{\tilde{\mathbf{d}}} f^{(k)}(\mathbf{x}^{(k)}) = \tilde{\mathbf{G}}^{(k)-T} \tilde{\nabla}_{\boldsymbol{\varphi}} h^{(k)}(\boldsymbol{\varphi}) \Big|_{\boldsymbol{\varphi}=0} \text{ where } h(\boldsymbol{\varphi}) = f(\mathbf{x}^{(k)} + \tilde{\mathbf{G}}^{(k)} \boldsymbol{\varphi}). \quad (5)$$

We start with  $\tilde{\mathbf{G}}^{(0)} = \mathbf{I}_n$ . This transformation (5) is a rotation/reflection using a different basis, and the length of the gradient is invariant due to the property of orthogonal matrices. The simplicity of this technique is represented by its easy implementation. The gradient is just computed in a different coordinate system using a simple reparamaterization of function  $f$  and a matrix-vector multiplication. We label this new approach by Smart Gradient.

**Simple example for illustration** Assume we have two iterations representing two  $\mathbf{x}$  positions of a two-dimensional function,

$$\mathbf{x}^{(0)} = (1.78, 2.82)^T, \mathbf{x}^{(1)} = (1.89, 4.62)^T, \mathbf{x}^{(2)} = (11.54, 4.15)^T, \dots$$

At each position  $k$ , we form  $\mathbf{G}^{(k)}$ . This matrix is initialized as the identity matrix, so at the first position  $k = 0$  the first direction used is  $\mathbf{d}_1^{(0)} = (1, 0)^T$  and the second  $\mathbf{d}_2^{(0)} = (0, 1)^T$ . When reaching the second position, we get a new direction  $\Delta \mathbf{x}^{(1)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)}$ , so matrix  $\mathbf{G}^{(1)}$  is constructed as follow:

$$\mathbf{G}^{(1)} = [\mathbf{x}^{(1)} - \mathbf{x}^{(0)} | \mathbf{d}_1^{(0)}] = \begin{pmatrix} 0.11 & 1 \\ 1.80 & 0 \end{pmatrix}$$

This matrix is orthonormalized using the MGS algorithm, such that the columns of  $\mathbf{G}^{(1)}$ ,  $\mathbf{g}_1^{(1)}$  and  $\mathbf{g}_2^{(1)}$ , are updated to form  $\tilde{\mathbf{G}}^{(1)}$  by columns  $\tilde{\mathbf{d}}_1^{(1)}$  and  $\tilde{\mathbf{d}}_2^{(1)}$ ,

$$\tilde{\mathbf{d}}_1^{(k)} = \frac{\Delta \mathbf{x}^{(k)}}{\|\Delta \mathbf{x}^{(k)}\|} \text{ and } \tilde{\mathbf{d}}_2^{(k)} = \frac{\Delta \mathbf{x}^{(k-1)} - (\mathbf{g}_1^{(k)} \cdot \Delta \mathbf{x}^{(k-1)}) \frac{\Delta \mathbf{x}^{(k-1)}}{\|\Delta \mathbf{x}^{(k-1)}\|}}{\left\| \Delta \mathbf{x}^{(k-1)} - (\mathbf{g}_1^{(k)} \cdot \Delta \mathbf{x}^{(k-1)}) \frac{\Delta \mathbf{x}^{(k-1)}}{\|\Delta \mathbf{x}^{(k-1)}\|} \right\|},$$

where  $\mathbf{g}_i^{(k)}$  is the  $i$ th column of  $\mathbf{G}^{(k)}$ , and

$$\tilde{\mathbf{G}}^{(1)} = \begin{pmatrix} 0.0601 & 0.9981 \\ 0.9981 & -0.0610 \end{pmatrix}.$$

At the next position,  $\mathbf{x}_2$ , the same procedure for the two updates are computed,

$$\mathbf{G}^{(2)} = [\mathbf{x}^{(2)} - \mathbf{x}^{(1)} | \mathbf{d}_1^{(1)}] = \begin{pmatrix} 9.65 & 0.0610 \\ -0.47 & 0.9981 \end{pmatrix}$$

then after MGS,

$$\tilde{\mathbf{G}}^{(2)} = \begin{pmatrix} 0.9989 & 0.0486 \\ -0.0486 & 0.9989 \end{pmatrix}$$

The columns of  $\tilde{\mathbf{G}}^{(2)}$  are orthogonal unit vectors, and they are used as the new set of axes in an orthogonal coordinate system in this unitary transformation of the gradient. For each formed  $\tilde{\mathbf{G}}^{(k)}$ , the gradient is estimated using (5). See Figure 2

for the plotted directions used to compute Vanilla Gradient and Smart Gradient (the first two iterations).

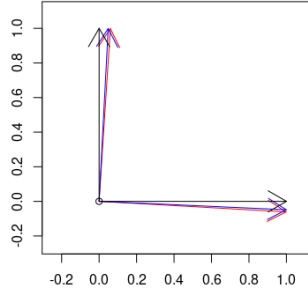


FIGURE 2. Directions used to compute Vanilla Gradient (black) - Smart Gradient (iteration 1 - red) and Smart Gradient (iteration 2 - blue).

### 3. Numerical Experiments.

**3.1. Test Functions.** In this section, some experiments are carried out to compare the performance of using the proposed approach. We chose to compare our approach with another numerical gradient estimator based on the first order central differences.

For each test function, we use `optim` function in `stats` package in R with the BFGS method [3] to get the optimum, and the average squared difference (MSE) between the exact gradient (calculated analytically) and the estimated gradient (using Vanilla Gradient and the Smart Gradient) is calculated at each iteration. The experiment is repeated 100 times with a random initial value.

One of the popular test problems for unconstrained optimization is the Extended Rosenbrock Function, which is unimodal, differentiable and has an optimum in a narrow parabolic valley, see Figure 1. Another one is the Extended Freudenstein Roth Function, and the equations for both functions are here:

1. Extended Rosenbrock Function:

$$f(\mathbf{x}) = \sum_{i=1}^{n/2} 100(x_{2i} - x_{2i-1}^2)^2 + (1 - x_{2i-1})^2 \quad (6)$$

2. Extended Roth Freudenstein Function:

$$f(\mathbf{x}) = \sum_{i=1}^{n/2} \left( -13 + x_{2i-1} + x_{2i}(x_{2i}(5 - x_{2i}) - 2) \right)^2 + \left( -29 + x_{2i-1} + x_{2i}(x_{2i}(x_{2i} + 1) - 14) \right)^2 \quad (7)$$

Comparison results are summarized in Table (1) showing an improvement (ratio of the errors of VG and SG) of 2.5 for dimension 5 and at least 3.5 for the other two higher dimensions in Rosenbrock function.

In Figure 3, the two estimated gradients started with almost the same MSE in the first some iterations, then gradually the MSE of the Smart Gradient decreases

$\mathbf{x}$ dimension	Average MSE Vanilla Gradient	Average MSE Smart Gradient	Improvement
<b>Extended Rosenbrock Function</b>			
5	2.60e-04	1.04e-04	2.5
10	2.71e-04	0.78e-04	3.47
25	2.80e-04	0.49e-04	5.71
<b>Extended Freudenstein Roth Function</b>			
5	2.26e-04	1.39e-04	1.63
10	2.78e-04	1.42e-04	1.96
25	2.84e-04	1.25e-04	2.27

TABLE 1. The average MSE when using Vanilla and Smart Gradient approaches compared to the exact gradient for different dimensional Rosenbrock and Roth functions

showing a clear improvement, and it stays lower than the MSE of the Vanilla Gradient till the end of the iterations. Matrix  $\mathbf{G}^{(k)}$  needs at least  $n$  iterations to be filled properly with informative directions, and improvements of the rounding off errors become evident after the number of iterations exceeds the dimension of  $\mathbf{x}$ . This matrix continues rolling up for different  $k$  until the optimum is found.

It is clear that the gradient is calculated more accurately by Smart Gradient for all the functions in this optimization framework, as the prior information we have about the the descent directions boosts this accuracy. We examined the use of Smart Gradients on range of different test functions, all showing similar behaviour and overall improvement to the examples reported.

**3.2. Smart Hessian.** The Smart Gradient technique can be extended and applied on Hessians. Assume for a continuous function  $f(\mathbf{x})$ , the second partial derivatives exit, we can estimate the Hessian  $\nabla^2 f^{(k)}(\mathbf{x}^{(k)})$  of this objection function at iteration  $k$  based on some directions  $\{\tilde{\mathbf{d}}_1^{(k)} | \tilde{\mathbf{d}}_2^{(k)} | \dots | \tilde{\mathbf{d}}_n^{(k)}\}$  as we did for gradient in formulas (4) and (5),

$$\tilde{\nabla}_{\tilde{\mathbf{d}}}^2 f^{(k)}(\mathbf{x}^{(k)}) = \tilde{\mathbf{G}}^{(k)-T} \tilde{\nabla}^2 h^{(k)}(\boldsymbol{\varphi}) \Big|_{\boldsymbol{\varphi}=\mathbf{0}} \tilde{\mathbf{G}}^{(k)T} \text{ where } h(\boldsymbol{\varphi}) = f(\mathbf{x}^{(k)} + \tilde{\mathbf{G}}^{(k)}\boldsymbol{\varphi}) \quad (8)$$

One important application to this Smart Hessian is computing the Hessian of a function at its mode. This adaptive technique can help describing better the curvature of a function at its optimum, using the last  $n$  descent directions.

## 4. Applications.

**4.1. Optimization of a Target Function.** In this example, we use Smart Gradient technique to get the optimum of different dimensional Extended Rosenbrock Function using `optim` function in `stats` package in R, and compare it with Vanilla Gradient. We compute the two gradient with different step-size, in the range of  $[10^{-8} - 10^{-1}]$ , using first order central difference, and we plot the results as the function of MSE of the converged optimum  $\mathbf{x}^*$  and the true optimum value  $\mathbf{1}$ , see Figure 4.

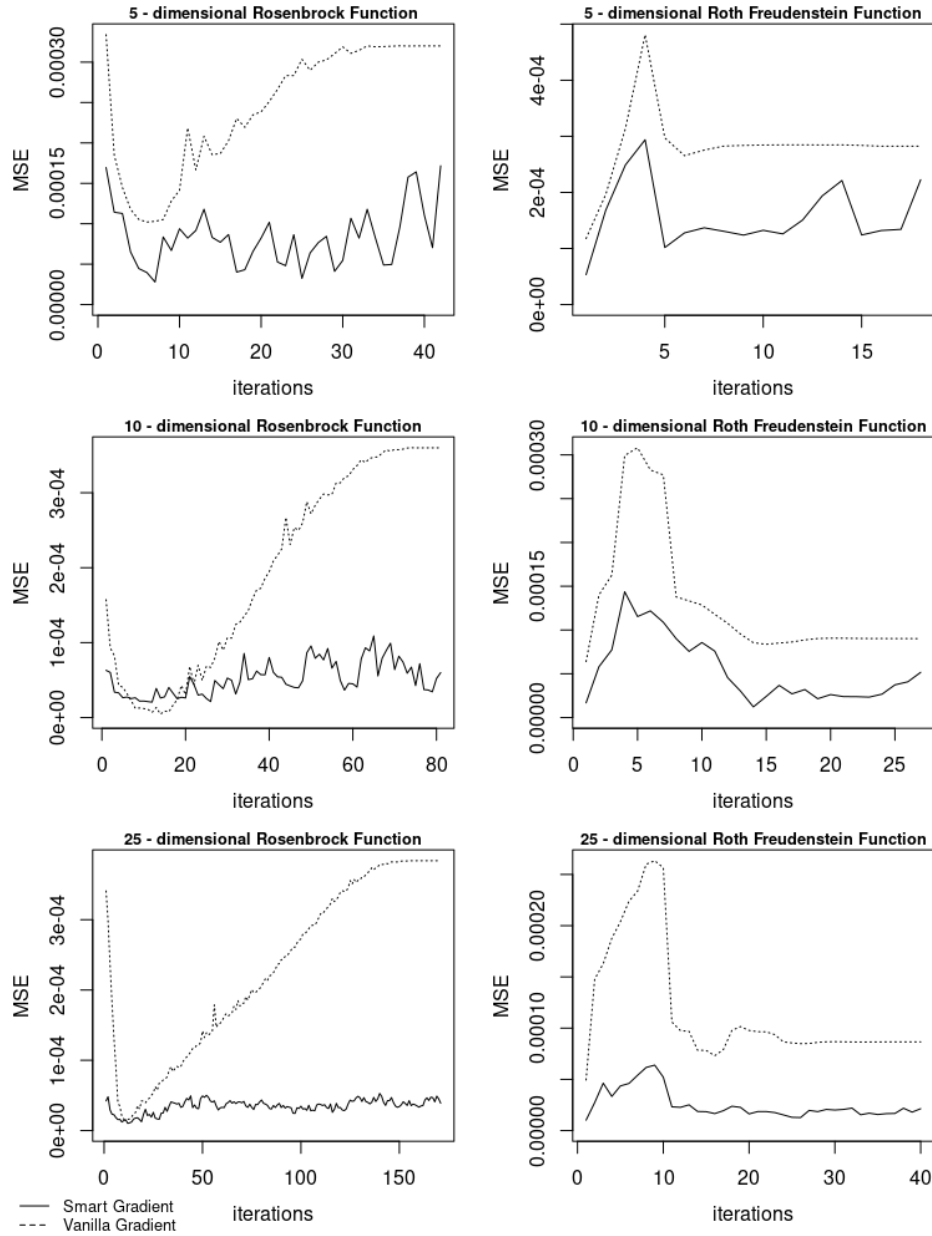


FIGURE 3. MSE for Smart and Vanilla Gradients at each iteration using Different Dimensional Rosenbrock and Roth Functions.

Also, we present in Figure 4 the computational time needed to reach the optimum when using Smart Gradient compared to Vanilla Gradient, as function of the step-size. Smart Gradient shows a higher computational cost compared to Vanilla gradient. However, this cost is still small compared to the accuracy gained, given that  $\mathbf{x}$  has moderate dimension. The MSE of the converged  $\mathbf{x}^*$  still has better overall accuracy.



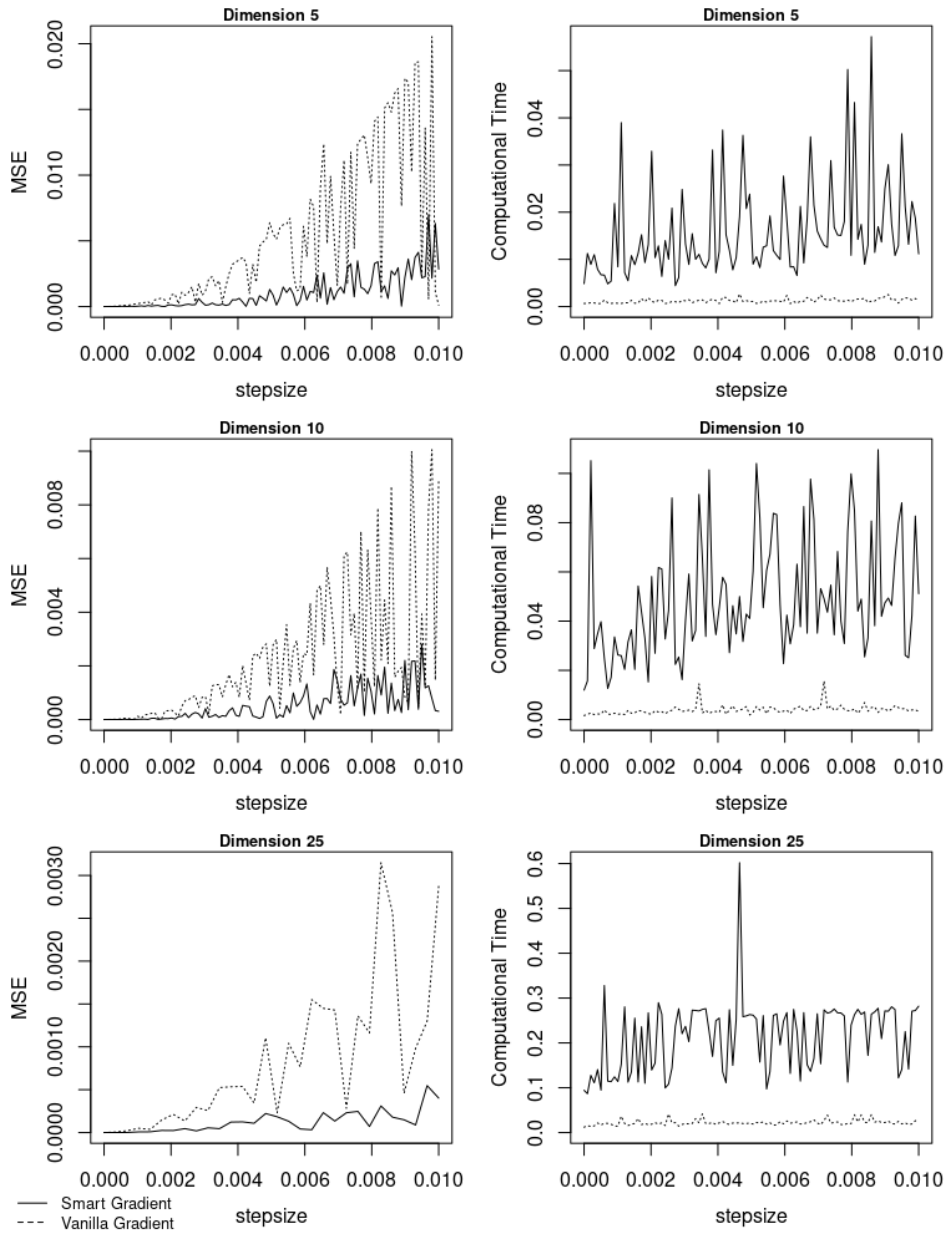


FIGURE 4. MSE and computational time (to reach optimum) as the function of step-size for different dimensional Rosenbrock Function.

**4.2. Hierarchical Optimization of a Hyperparameter.** A motivating application of the Smart Gradient technique is approximate Bayesian inference of a latent Gaussian model using Integrated Nested Laplace Approximation (INLA) [8], [9]. INLA uses combinations of analytical approximations and numerical integration to obtain approximate posterior distributions of the parameters.

4.2.1. *Brief overview of INLA.* Latent Gaussian models (LGMs) form a class of regression models that include generalized linear models (GLM) like Poisson regression, generalized additive mixed models (GAMM), temporal and spatial models, amongst many others. An LGM is defined for data  $\mathbf{y}$ , a Gaussian latent field  $\mathbf{z}$  and hyperparameters  $\mathbf{x}$ ,

$$\begin{aligned} \mathbf{y}|\mathbf{z}, \mathbf{x} &\sim \pi(h(\mathbf{z}), \mathbf{x}) \\ \mathbf{z}|\mathbf{x} &\sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1}(\mathbf{x})) \\ \mathbf{x} &\sim \pi(\mathbf{x}), \end{aligned} \tag{9}$$

where  $\pi(\cdot, \cdot), \pi(\cdot)$  are density functions,  $h$  is a link function, and  $\mathbf{Q}(\mathbf{x})$  is the precision matrix of the latent field  $\mathbf{z}$ .

The main goal of INLA is to approximate the posteriors of  $\mathbf{z}$  and  $\mathbf{x}$ , i.e.  $\pi(\mathbf{z}|\mathbf{y})$  and  $\pi(\mathbf{x}|\mathbf{y})$ , respectively. INLA performs these approximations in two steps - firstly, optimize  $\pi(\mathbf{x}|\mathbf{y})$  and then use the optimum value of  $\mathbf{x} = \mathbf{x}^*$  to optimize  $\pi(\mathbf{z}|\mathbf{y}, \mathbf{x})$  to find  $\mathbf{z}^*$ , and subsequently approximate  $\pi(\mathbf{z}|\mathbf{y})$  with numerical integration. Note that the dimension of  $\mathbf{z}$  is often large ( $10^3 - 10^5$ ) but that of  $\mathbf{x}$  is often small ( $0 - 20$ ). The BFGS method [3] is used to estimate the optimum value of  $\mathbf{x}$ . The iterative nature of INLA and the size of  $\mathbf{z}$  prohibits exact function evaluations.

Assume  $\mathbf{x}$  has a prior distribution  $\pi(\mathbf{x})$ , we use the marginal posterior of  $\mathbf{x}$  to find the mode  $\mathbf{x}^*$ ,

$$\pi(\mathbf{x}|\mathbf{y}) = \frac{\pi(\mathbf{z}, \mathbf{x}|\mathbf{y})}{\pi(\mathbf{z}|\mathbf{y}, \mathbf{x})} = \frac{\pi(\mathbf{x})\pi(\mathbf{z}|\mathbf{x})\pi(\mathbf{y}|\mathbf{z})}{\pi(\mathbf{z}|\mathbf{y}, \mathbf{x})}.$$

INLA method constructs the approximate  $\tilde{\pi}(\mathbf{x}|\mathbf{y})$  by replacing the full conditional of  $\mathbf{z}$ :  $\pi(\mathbf{z}|\mathbf{y}, \mathbf{x})$  by its Gaussian approximation, then we write the Laplace approximation of  $\pi(\mathbf{x}|\mathbf{y})$ , and evaluate it at the mode  $\mathbf{z}^*$ ,

$$\mathbf{x}^* \approx \underset{\mathbf{x}}{\operatorname{argmin}} - \left( \log \pi(\mathbf{y}|\mathbf{z}^*, \mathbf{x}) + \log \pi(\mathbf{z}^*|\mathbf{x}) + \log \pi(\mathbf{x}) - \log \pi(\mathbf{z}^*|\mathbf{y}, \mathbf{x}) \right). \tag{10}$$

The gradient of the objective function in INLA applications is not exact and needs to be estimated numerically. A cheap and efficient technique is Smart Gradient, that shows a great improvement in the MSE with small added cost. Assume for example the dimension of the hyperparameter  $\mathbf{x}$  is 5 and that of the latent field is  $10^3$ . Adding  $5^2$  and  $5^3$  operations to estimate the gradient at every iteration has essentially small cost compared to (1) the gained accuracy of this gradient and (2) the high number of operations needed at every iteration for the inner optimization to get  $\mathbf{z}^*$ .

4.2.2. *Probit regression example.* Probit regression is used when we have binary data and the goal is to model the probability of a success (a value of 1). To illustrate the Smart Gradient technique in a complex optimization problem we consider a continuously indexed spatial probit regression model,

$$\begin{aligned} y_i|\mathbf{p} &\sim \operatorname{Bin}(p_i) \\ \operatorname{logit}(p_i)|\beta_0, \beta_1, \mathbf{w} &= \beta_0 + \beta_1 v_i + \mathbf{A}\mathbf{w}, \end{aligned} \tag{11}$$

where  $\mathbf{p}$  is the vector of probabilities,  $\mathbf{v}$  is a vector of covariates,  $\mathbf{w}$  is a Gaussian spatial random field and  $\mathbf{A}$  is projector matrix that links the spatial field with the data. The fixed effects  $\beta_0$  and  $\beta_1$  follow a weakly informative Gaussian distribution a priori, with constant precision  $1e^{-3}$ . The Gaussian spatial field is assumed to have

a Matérn correlation structure with smoothness parameter  $\nu$ , scaling parameter  $\kappa$  and marginal variance  $\sigma_w^2$  such that  $\mathbf{w} \sim N(\mathbf{0}, \mathbf{R})$ , with the entries of  $\mathbf{R}$  as,

$$R_{ij} = \frac{\sigma_w^2}{2^{\nu-1}\Gamma(\nu)} (\kappa \|\mathbf{s}_i - \mathbf{s}_j\|)^\nu K_\nu(\kappa \|\mathbf{s}_i - \mathbf{s}_j\|) \quad (12)$$

for locations  $\mathbf{s}_i$  and  $\mathbf{s}_j$ , where  $K_\nu(\cdot)$  is the modified Bessel function of the second kind and  $\|\cdot\|$  is the Euclidean distance.

The continuously indexed Matérn spatial field  $\mathbf{w}$  is defined as the weak solution of the following Stochastic Partial Differential Equation (SPDE) [4],

$$(\kappa^2 - \Delta)^{\frac{\alpha}{2}} \mathbf{w}(\mathbf{s}) = \mathbf{W}(\mathbf{s}), \quad (13)$$

with  $\alpha = \nu + 1$  and where  $\mathbf{W}(\mathbf{s})$  is a white noise Gaussian stochastic process in space (i.e.  $\mathbb{R}^2$ ) with unit variance, and  $\Delta$  is the Laplacian operator. The parameters are reparameterized for orthogonality and interpretability to a range parameter  $\rho = \sqrt{8\nu}/\kappa$  and marginal variance  $\sigma_w^2 = \Gamma(\nu)/\Gamma(\nu + 1)4\pi\kappa^{2\nu}$ .

The solution  $\mathbf{w}(\mathbf{s})$  to (13) is found through the finite element method by discretizing the continuous space into triangles (which forms a mesh) and mapping the observed locations with a projector matrix  $\mathbf{A}$ , to this mesh composed of triangles.

For  $m$  triangles in the mesh, the latent field is  $\mathbf{z} = \{\beta_0, \beta_1, \mathbf{w}\}$  and of dimension  $m + 2$ . We use the unit square as a spatial domain and the accompanying mesh is illustrated in Figure 5. Note that we extend the mesh beyond the study area to avoid any boundary effects. Independent Gaussian priors with precisions  $\tau_{\beta_0}$  and  $\tau_{\beta_1}$  for  $\beta_0$  and  $\beta_1$ , respectively, results in a latent Gaussian prior with precision matrix

$$\mathbf{Q}(\mathbf{x}) = \begin{pmatrix} \tau_{\beta_0} & 0 & 0 \\ 0 & \tau_{\beta_1} & 0 \\ 0 & 0 & \mathbf{R}^{-1}(\mathbf{x}) \end{pmatrix},$$

where  $\mathbf{x} = (\log \sigma_w, \log \rho)$ .

We simulate a sample of size 1000, with  $\beta_0 = 1/2, \beta_1 = 1, \sigma_w^2 = 0.5, \nu = 0.5$  and  $\kappa = 10$  and represent the observed locations in Figure 5. Now we perform Bayesian inference for this spatial probit regression model using the `inla()` function. In this inferential procedure, we need to find the modes of the hyperparameters  $\kappa, \tau_w = \sigma_w^{-2}$ , so that the marginals of the elements of the Gaussian latent field  $\beta_0$  and  $\beta_1$  can be estimated. We use (10) to find  $\mathbf{x}^*$ .

When the data is generated, we specify the true  $\mathbf{x}$  to be  $(-\log \sqrt{2}, -\log 5)$  and we use the BFGS method with first order central difference to reach the maximum. We use INLA package in R to fit the model using Vanilla and Smart gradients. The absolute mean error is calculated in each case, and the number of calls to reach  $\boldsymbol{\theta}^*$  is stored. We repeated the experiment 1000 times, for Vanilla Gradient the average absolute mean error is 0.26224 and the average number of calls is 220, whereas for Smart Gradient the average error is 0.26110 and average calls is 206. When using Smart Gradient, the optimum is reached with less number of calls of the objective function and with an improvement in the optimum value.

**5. Discussion and further considerations.** Gradients are an important tool that is used in various applied mathematical and statistical methods. This wide use of gradients necessitates the search for techniques that improve its estimation.

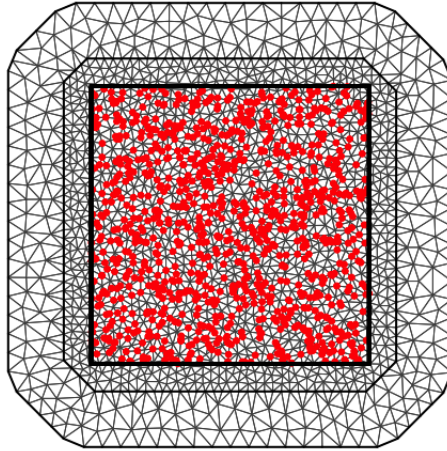


FIGURE 5. Triangulated mesh with 1000 locations on the unit square

Here, we presented a simple framework that enhances the accuracy of gradient estimation within optimization context using the most recent differences between  $\mathbf{x}$  positions as directions to compute the gradient.

In general, using more precise gradient will lead to a better and more robust performance for gradient based optimization methods. For a moderate number of dimension, say  $\leq 25$ , Smart Gradient method shows improvement in its accuracy with essentially small cost. One useful application for the presented approach is the Hierarchical optimization. That is the case when evaluating the function at two points or more to compute the gradient is costly because of the inner optimization, and using Smart Gradient with its added small cost to gain more accuracy is worth it.

The proposed method is naturally extended to improve the estimate of the Hessian through the attained descent directions. It is implemented in the accompanying R package `smartGrad` available on the github at [esmail-abdulfattah/Smart-Gradient](https://github.com/esmail-abdulfattah/Smart-Gradient). Additionally, `smartGrad` can be used to enhance a user-defined gradient formula through the function `makeSmart`, see Appendix A. C++ code is also available on github.

**References**

- [1] Besag, J. (1975). Statistical analysis of non-lattice data. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 24(3):179–195.
- [2] Depner, J. S. and Rasmussen, T. C. (2016). *Hydrodynamics of Time-periodic Groundwater Flow: Diffusion Waves in Porous Media*, volume 224. John Wiley & Sons.
- [3] Fletcher, R. (1988). Practical methods of optimization.
- [4] Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between gaussian fields and gaussian markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(4):423–498.
- [5] Nocedal and dan Stephen J. Wright, J. (2020). Numerical optimization, 2nd edition.
- [6] Nocedal, J. and Wright, S. J. (1999). Numerical optimization.
- [7] Picheny, V., Wagner, T., and Ginsbourger, D. (2013). A benchmark of kriging-based infill criteria for noisy optimization. *Structural and Multidisciplinary Optimization*, 48:607–626.
- [8] Rue, H., Martino, S., and Chopin, N. (2009). Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of The Royal Statistical Society Series B-statistical Methodology*, 71:319–392.
- [9] Rue, H., Riebler, A., Sørbye, S., Illian, J., Simpson, D. P., and Lindgren, F. (2016). Bayesian computing with inla: A review.
- [10] Thomas, G., Weir, M. D., Hass, J., and Giordano, F. (2005). Thomas’ calculus early transcendentials (11th edition) (thomas series).

## Appendix A. smartGrad Installation and Examples.

### A.1. Installation.

```
library("devtools")
install_github("esmail-abdulfattah/Smart-Gradient",
              subdir = "smartGrad")
```

A.2. **makeSmart Function.** To illustrate how to use this function, we use the Extended Rosenbrock function of dimension  $n$ . It has global minimum at  $\mathbf{x}^* = \mathbf{1}$ . We estimate the gradient of this objective function  $f$  using a simple central difference method, with step size  $10^{-3}$ .

```
myfun <- function(x) {
  res <- 0.0
  for(i in 1:(length(x)-1))
    res <- res + 100*(x[i+1] - x[i]^2)^2 + (1-x[i])^2
  return(res)
}

mygrad <- function(fun,x){
  h = 1e-3
  grad <- numeric(length(x))
  for(i in 1:length(x)){
    e = numeric(length(x))
    e[i] = 1
    grad[i] <- (fun(x+h*e) - fun(x-h*e))/(2*h)
  }
  return(grad)
}
```

A user can change a numerical gradient function `mygrad` to a SMART numerical gradient function `mySmartgrad`, and then this SMART function can be used instead. We use the `optim` function from `stats` package with BFGS algorithm.

```
library("stats")
library("smartGrad")
x_dimension = 5
x_initial = rnorm(x_dimension)
result <- optim(par = x_initial,
               fn = myfun,
               gr = makeSmart(fn = myfun,gr = mygrad),
               method = c("BFGS"))
```

Received xxxx 20xx; revised xxxx 20xx.

*E-mail address:* [esmail.abdulfattah@kaust.edu.sa](mailto:esmail.abdulfattah@kaust.edu.sa)

*E-mail address:* [janet.vannierkerk@kaust.edu.sa](mailto:janet.vannierkerk@kaust.edu.sa)

*E-mail address:* [haavard.rue@kaust.edu.sa](mailto:haavard.rue@kaust.edu.sa)