
SMART GRADIENT - AN ADAPTIVE TECHNIQUE FOR IMPROVING GRADIENT ESTIMATION

A PREPRINT

✉ **Esmail H. Abdul Fattah***

Department of Statistics
King Abdullah University of Science and Technology
Thuwal, 23955, Makkah
esmail.abdulfattah@kaust.edu.sa

Janet Van Niekerk

Department of Statistics
King Abdullah University of Science and Technology
Thuwal, 23955, Makkah
janet.vanniekerk@kaust.edu.sa

Håvard Rue

Department of Statistics
King Abdullah University of Science and Technology
Thuwal, 23955, Makkah
haavard.rue@kaust.edu.sa

June 15, 2021

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 require gradients are optimization techniques such as stochastic gradient descent, Newton’s method and trust region methods. However, these methods usually requires 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++.

Keywords Adaptive Technique · Gradient Estimation · Numerical Gradient · Optimization · Vanilla Gradient Descent

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} [1]. 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 [2] 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,

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

*Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

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

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

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 [3] 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^{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 α_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 [4]:

1. Gradient Descent: $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$.
2. Newton's Method: $\mathbf{d}^{(k)} = -\nabla^2 f(\mathbf{x}^{(k)})^{-1} \nabla f(\mathbf{x}^{(k)})$.
3. 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 $\mathbf{x}^{(k)}$, $\nabla^2 f(\mathbf{x}^{(k)})$.

The choice of the step size $\{\alpha^{(k)}\}$ is obtained using either exact or inexact line search [1].

We summarize and relate four popular gradient computational frameworks next [5].

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

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

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,

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

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

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

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 [6] $f(\mathbf{x}) = (1 - x_1)^2 + 100 \times (x_2 - x_1^2)^2$, see Figure 2, using various bases. We use the central difference method of first order of step length 10^{-3} to get the gradient estimate at $\mathbf{x}^{(0)} = (-0.29, 0.40)^T$. 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.

From Figure 1 we observe opposite trends between the error in the gradient estimates and the magnitude of the gradient itself: the MSE error in the gradient estimate is low when the magnitude of the gradient is high. Although such a claim does not hold in general, we found empirically it holds “almost everywhere”, 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} .

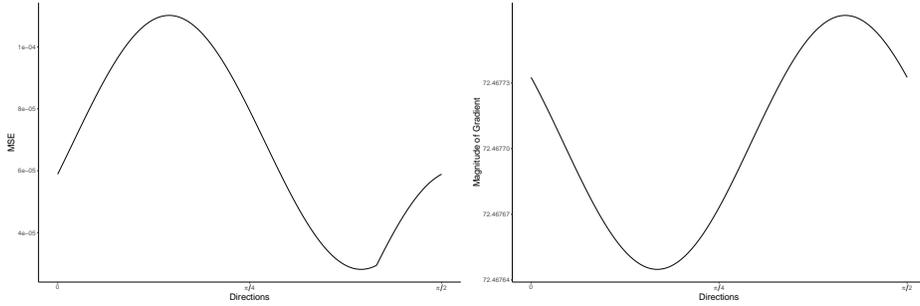


Figure 1: Estimating Gradient using different Directions

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 some applications within the R-INLA package. The paper is concluded by a discussion in Section 4.

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,

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

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),

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

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 reparameterization 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 of resulting in the set of \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 $\tilde{\mathbf{d}}_2^{(1)}$,

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

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

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

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).

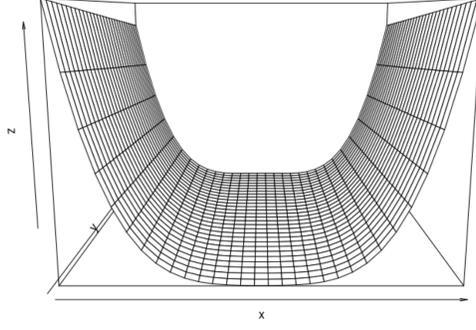


Figure 2: Rosenbrock Function

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 [2] 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 2. Another one is the Extended Freudenstein Roth Function, and the equations for both functions are here:

1. Extended Rosenbrock Function:

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

2. Extended Roth Freudenstein Function:

$$f(\mathbf{x}) = \sum_{i=1}^{n/2} \left(-13 + x_{2i-1} + x_{2i}(x_{2i}(5 - x_{2i}) - 2) \right)^2 + \left(-29 + x_{2i-1} + x_{2i}(x_{2i}(x_{2i} + 1) - 14) \right)^2 \quad (6)$$

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 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. For higher order finite different schemes with increased accuracy, the Smart Gradients had similar overall MSE than estimating Vanilla gradients directly. This is reasonable, as Smart Gradients does not offer any improvement in the limit.

x dimension	Average MSE Vanilla Gradient	Average MSE Smart Gradient	Improvement
Extended Rosenbrock Function			
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
Extended Freudenstein Roth Function			
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

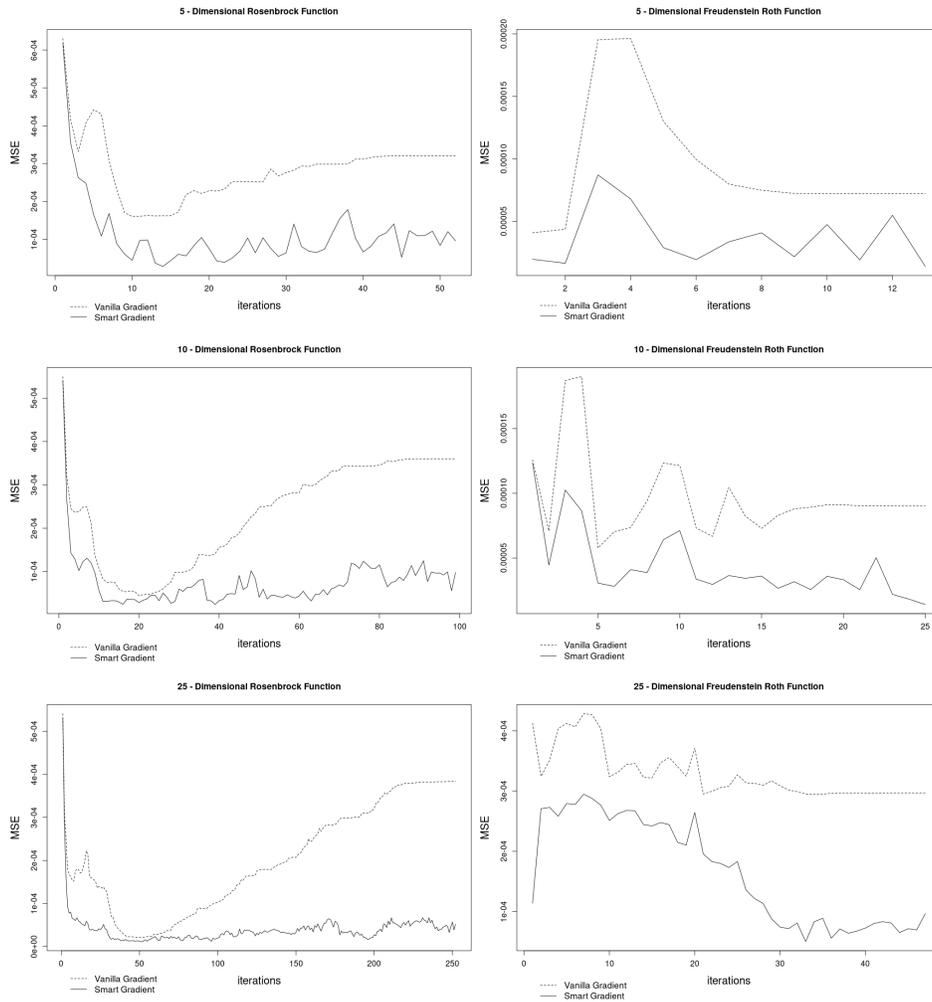


Figure 3: MSE for Smart and Vanilla Gradients at each iteration using Different Dimensional Rosenbrock and Roth Functions

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 exist, 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),

$$\tilde{\nabla}_{\mathbf{d}}^2 f^{(k)}(\mathbf{x}^{(k)}) = \tilde{\mathbf{G}}^{(k)-T} \tilde{\nabla}^2 h^{(k)}(\boldsymbol{\varphi}) \Big|_{\boldsymbol{\varphi}=\mathbf{0}} \tilde{\mathbf{G}}^{(k)T} \text{ where } h(\boldsymbol{\varphi}) = f(\mathbf{x}^{(k)} + \tilde{\mathbf{G}}^{(k)}\boldsymbol{\varphi}) \quad (7)$$

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.

3.3 Autoregressive time-series model with R-INLA

A motivation application to the Smart Gradient technique is fitting a model in Bayesian Inference using Integrated Nested Laplace Approximation (INLA) method [8], [9]. INLA uses combinations of analytical approximations and numerical integration to obtain approximated posterior distributions of the parameters. It uses the BFGS method [2] to reach the optimum value of the hyperparameter vector $\boldsymbol{\theta}$ in the model. At each iteration, an inner optimization takes place to approximate the posterior distribution of the latent field \mathbf{x} by Gaussian approximation and get the mode \mathbf{x}^* . This prohibits exact function evaluations.

The proposed technique of Smart Gradient is already used as the default option in INLA package in R to optimize the hyperparameters, using the following argument `inla(control.inla=list(use.directions = TRUE))` in `inla` function. We show in the next example how this method is used to fit a time series model using the INLA method.

Consider the R dataset of monthly totals of international airlines passengers between 1949 to 1960, see Figure 4, one possible option is to fit the response as autoregressive model of order 1,

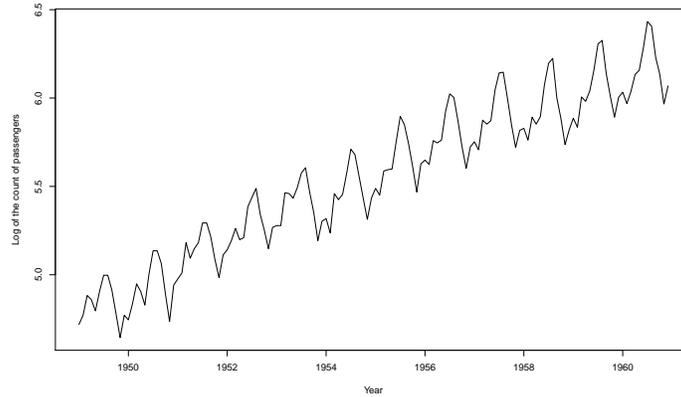


Figure 4: The log of air passengers from 1949 to 1961

$$\log(\mathbf{y}) \sim \beta \mathbf{t} + \mathbf{u} + \boldsymbol{\epsilon}$$

$$u_i \sim \phi u_{i-1} + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(0, \tau_u^{-1}), i = 2, \dots, n, \text{ and} \\ u_1 \sim \mathcal{N}(0, (\tau_u(1 - \phi^2))^{-1})$$

where \mathbf{y} is the number of air passengers, \mathbf{t} is the year, ε_j follows skew-normal distribution with zero mean, $j = 1, \dots, n$, standardised skewness γ and τ_y as the precision parameter. The fixed effect β on \mathbf{t} follow a weakly informative Gaussian distribution a priori, with constant τ_β as precision. The four hyperparameters $(\tau_y, \gamma, \tau_u, \phi)$ in our model are transformed to $\boldsymbol{\theta} = (\theta_1 = \log(\tau_y), \theta_2 = \log((1 + \gamma/0.988)/(1 - \gamma/0.988)), \theta_3 = \log(\tau_u(1 - \phi^2)), \theta_4 = \log((1 + \phi)/(1 - \phi)))$ for unconstrained optimization.

In this inferential procedure we need to find the mode of the hyperparameters $\boldsymbol{\theta}$, so that the marginals of the elements of the Gaussian latent field $\mathbf{x}^T = (\beta, \mathbf{u}^T)$ can be estimated. The objective function $f(\boldsymbol{\theta}|\mathbf{y})$,

$$\log f(\boldsymbol{\theta}|\mathbf{y}) \approx -\log \pi(\mathbf{y}|\mathbf{x}^*, \boldsymbol{\theta}) + \log \pi(\mathbf{x}^*|\boldsymbol{\theta}) + \log \pi(\boldsymbol{\theta}) - \pi(\mathbf{x}^*|\mathbf{y}, \boldsymbol{\theta})$$

where \mathbf{x}^* is the mode we get from the inner optimization of $\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$ as it is approximated by Gaussian. Using PC priors [10], we take the prior distributions for $\boldsymbol{\theta}$,

$$\begin{aligned} \log \pi(\boldsymbol{\theta}) &\approx -\log(0.2)e^{-\theta_1/2} - \theta_1/2 \\ &\quad + \log(\phi(\theta_2)) + \log(\Phi(10^{1/3}\theta_2)) \\ &\quad + \theta_3 - (5 \times 10^{-5})e^{\theta_3} - 0.075\theta_4^2 \end{aligned} \quad (8)$$

where $\phi(\cdot)$ is the standard normal distribution and $\Phi(\cdot)$ is the standard cumulative Gaussian distribution. With $\boldsymbol{\eta}$ as the linear predictor and $\mathbf{Q}^* = \mathbf{Q}(\boldsymbol{\theta}) + \mathbf{Q}_l$, where $\mathbf{Q}(\boldsymbol{\theta})$ is a combination of the AR1 model precision matrix $\mathbf{R}(\boldsymbol{\theta})$ and τ_β ,

$$\mathbf{Q}(\boldsymbol{\theta}) = \begin{pmatrix} \tau_\beta & 0 \\ 0 & \mathbf{R}(\boldsymbol{\theta}) \end{pmatrix},$$

and \mathbf{Q}_l is the precision we get from the model, then

$$\begin{aligned} \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) &\approx \exp\left(-\frac{\tau_\epsilon}{2}(\mathbf{y} - \boldsymbol{\eta})^T(\mathbf{y} - \boldsymbol{\eta})\right) \\ \pi(\mathbf{x}|\boldsymbol{\theta}) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T\mathbf{Q}(\boldsymbol{\theta})\mathbf{x}\right) \\ \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) &\approx \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T\mathbf{Q}^*(\mathbf{x} - \mathbf{x}^*)\right) \end{aligned}$$

Due to the intractable form of the objective function $f(\boldsymbol{\theta}|y)$, it is hard to get the exact gradient, so it needs to be estimated numerically. Here we use Smart Gradient integrated with the central difference method which is more accurate compared to using Vanilla Gradient with canonical basis.

Using Smart Gradient technique in this unconstrained optimization, we get $\boldsymbol{\theta}^* = (2454.55, -0.001, 43.483, 0.744)$, and the Smart Hessian of the objective function f at $\boldsymbol{\theta}^*$ can be calculated easily the same way.

3.4 Continuous spatial statistical model with the stochastic partial differential equation (SPDE) method

R-INLA

Consider the R dataset Leuk that features the survival times of patients with acute myeloid leukemia (AML) in Northwest England between 1982 to 1998. Exact residential locations and districts of the patients are known and indicated by the dots in Figure 5. The aim is to model the survival time based on various covariates \mathbf{X} and space \mathbf{s} , with

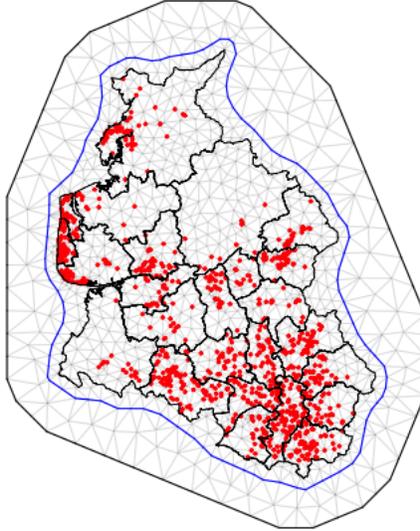


Figure 5: Exact residential locations of patients with AML

a Weibull model with shape α ,

$$t(\mathbf{s}) \sim \text{Weibull}(\lambda = \exp[\boldsymbol{\beta}\mathbf{X} + \mathbf{u}(\mathbf{s})], \alpha).$$

The shape parameter of the likelihood is also considered a hyperparameter (different to Section 3.3). To include a spatial element we use a Gaussian random effect \mathbf{u} with a Matern covariance structure that has marginal variance σ_u^2 and nominal range $r = 2/\kappa$ such that the entry of the precision matrix of \mathbf{u} between locations s_i and s_j is $\mathbf{R}_{ij} = \sigma_u^2 e^{-\kappa(s_i - s_j)^T (s_i - s_j)}$.

We use the finite element method and the mesh presented in Figure 5 to estimate this model (for more details regarding the SPDE approach to continuous spatial modeling see [11]). In this example we have incomplete observations in the sense that some patients are still alive at the end of the study and some times are thus censored, indicated in the variable \mathbf{c} . The data is thus $\mathbf{y} = \{\mathbf{t}, \mathbf{c}\}$ instead of a univariate observation as in Section 3.3.

Thus we have three hyperparameters in this model $\boldsymbol{\theta} = \{10 \log(\alpha), \log(2/\kappa), \log(\sigma_u^2)\}$, one from the likelihood and two from the spatial field which we need to optimize in order to obtain the marginal posteriors of the latent field. The objective function for $\boldsymbol{\theta}$ is obtained similar to that in Section 3.3 and again is intractable, necessitating a numerical gradient. We deploy Smart Gradient technique again and we calculate the optimal values for $\boldsymbol{\theta}$ as $\boldsymbol{\theta}^* = (-5.2246, 0.6530, 0.1599)$.

4 Discussion and further considerations

Gradients are important tool in various applied mathematical and statistical methods. This wide use of gradients necessitates the search for techniques that improve its estimation. 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 performance for gradient based optimization methods. For a moderate number of dimension, Smart Gradient method shows improvement in its accuracy with essentially no cost. For instance, a Smart Gradient technique with first order central difference method, where function is evaluated at only two evaluation points can perform as well as higher order differences methods which are computationally more costly.

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`. 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] J. Nocedal, S. J. Wright, Numerical optimization, 1999.
- [2] R. Fletcher, Practical methods of optimization, 1988.
- [3] G. Thomas, M. D. Weir, J. Hass, F. Giordano, Thomas' calculus early transcendentals (11th edition) (thomas series), 2005.
- [4] Nocedal, J. dan Stephen J. Wright, Numerical optimization, 2nd edition, 2020.
- [5] J. S. Depner, T. C. Rasmussen, Hydrodynamics of Time-periodic Groundwater Flow: Diffusion Waves in Porous Media, Vol. 224, John Wiley & Sons, 2016.
- [6] J. Besag, Statistical analysis of non-lattice data, Journal of the Royal Statistical Society: Series D (The Statistician) 24 (3) (1975) 179–195.
- [7] V. Picheny, T. Wagner, D. Ginsbourger, A benchmark of kriging-based infill criteria for noisy optimization, Structural and Multidisciplinary Optimization 48 (2013) 607–626.
- [8] H. Rue, S. Martino, N. Chopin, Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations, Journal of The Royal Statistical Society Series B-statistical Methodology 71 (2009) 319–392.
- [9] H. Rue, A. Riebler, S. Sørbye, J. Illian, D. P. Simpson, F. Lindgren, Bayesian computing with inla: A review, 2016.
- [10] D. P. Simpson, H. Rue, T. G. Martins, A. Riebler, S. Sørbye, Penalising model component complexity: A principled, practical approach to constructing priors, Statistical Science 32 (2014) 1–28.
- [11] E. T. Krainski, V. Gómez-Rubio, H. Bakka, A. Lenzi, D. Castro-Camilo, D. Simpson, F. Lindgren, H. Rue, Advanced spatial modeling with stochastic partial differential equations using R and INLA, CRC Press, 2018.

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"))
```