# Chapter 13
# Gradient Sampling Methods

One of the newest approaches in general NSO is to use gradient sampling algorithms developed by Burke et al. [51, 52]. The *gradient sampling method* (GS) is a method for minimizing an objective function that is locally Lipschitz continuous and smooth on an open dense subset $D \subset \mathbb{R}^n$. The objective may be nonsmooth and/or nonconvex. The GS may be considered as a stabilized steepest descent algorithm. The central idea behind these techniques is to approximate the subdifferential of the objective function through random sampling of gradients near the current iteration point. The ongoing progress in the development of gradient sampling algorithms (see e.g. [67]) suggests that they may have potential to rival bundle methods in the terms of theoretical might and practical performance. However, here we introduce only the original GS [51, 52].

## 13.1 Gradient Sampling Method

Let $f$ be a locally Lipschitz continuous function on $\mathbb{R}^n$, and suppose that $f$ is smooth on an open dense subset $D \subset \mathbb{R}^n$. In addition, assume that there exists a point $\bar{x}$ such that the level set $\text{lev}_{f(\bar{x})} = \{x \mid f(x) \leq f(\bar{x})\}$ is compact.

At a given iterate $x_k$ the gradient of the objective function is computed on a set of randomly generated nearby points $u_{kj}$ with $j \in \{1, 2, \ldots, m\}$ and $m > n + 1$. This information is utilized to construct a search direction as a vector in the convex hull of these gradients with the shortest norm. A standard line search is then used to obtain a point with lower objective function value. The stabilization of the method is controlled by the *sampling radius* $\varepsilon_k$ used to sample the gradients.

The pseudo-code of the GS is the following:

```
PROGRAM GS
   INITIALIZE  x₀ ∈ lev_f(x̄) ∩D,  ε₀ > 0,  m > n + 1,  ν₀ ≥ 0,  θ, μ ∈ (0, 1]
      and  α, β ∈ (0, 1);
   Set  k = 0;
   WHILE the termination condition is not met
      GRADIENT SAMPLING
         Sample  u_k1,...u_km  from  B̄(x; 1);
         Set  x_k0 = x_k and  x_kj = x_k + ε_k u_kj  for  j = 1,..., m;
         IF  x_kj ∉ D for some  j STOP;
         Set  G_k = {∇f(x_k1), ∇f(x_k2),..., ∇f(x_km)};
      END GRADIENT SAMPLING
      Compute  g_k = argmin_{g∈G_k} ‖g‖²;
      IF  ν_k = ‖g_k‖ = 0 STOP with the final solution  x_k;
      IF  ‖g_k‖ > ν_k  THEN
         Set  ν_{k+1} = ν_k and  ε_{k+1} = ε_k;
         Compute the search direction  d_k = −g_k/‖g_k‖;
         Find the step size  t_k = max α^p  such that
            f(x_k + α^p d_k) < f(x_k) − βα^p ‖g_k‖  and  p ∈ {1, 2,...};
      ELSE
         Set  t_k = 0,  ν_{k+1} = θν_k,  and  ε_{k+1} = με_k;
      END IF
      IF  x_k + t_k d_k ∈ D  THEN           Set  x_{k+1} = x_k + t_k d_k;
      ELSE
         Let  x̂^k  be any point on  B̄(x; ε_k)  satisfying  x̂_k + t_k d_k ∈ D
            and  f(x̂_k + t_k d_k) < f(x̂_k) − βt_k ‖g_k‖  (such a point exists
            due to continuity of  f);
         Set  x_{k+1} = x̂_k + t_k d_k;
      END IF
      Set  k = k + 1;
   END WHILE
   RETURN final solution  x_k;
END PROGRAM GS
```

Note that the probability to obtain a point $x_{kj} \notin D$ is zero in the above algorithm. In addition, it is reported in [52] that it is highly unlikely to have $x_k + t_k d_k \notin D$.

The GS algorithm may be applied to any function $f : \mathbb{R}^n \to \mathbb{R}$ that is continuous on $\mathbb{R}^n$ and differentiable almost everywhere. Furthermore, it has been shown that when $f$ is locally Lipschitz continuous, smooth on an open dense subset $D$ of $\mathbb{R}^n$, and has bounded level sets, the cluster point $\bar{x}$ of the sequence generated by the GS with fixed $\varepsilon$ is $\varepsilon$-stationary with probability 1 (that is, $\mathbf{0} \in \partial_\varepsilon^G f(\bar{x})$, see also Definition 3.3 in Part I). In addition, if $f$ has a unique $\varepsilon$-stationary point $\bar{x}$, then the set of all cluster points generated by the algorithm converges to $\bar{x}$ as $\varepsilon$ is reduced to zero.