

Optimization for Data Science

Lec 14: Randomized Smoothing

Yaoliang Yu



UNIVERSITY OF
WATERLOO

FACULTY OF MATHEMATICS
**DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE**

Problem

Constrained minimization:

$$\min_{\mathbf{w} \in C \subseteq \mathbb{R}^d} f(\mathbf{w})$$

- C is closed convex and f is (non)convex
- Can only evaluate the function value $f(\mathbf{w})$ but not the (sub)gradient
- Zero-th order method (a.k.a. gradient-free or derivative-free)
- For most (if not all) functions in practice, computing the function value (a scalar) costs as much as computing a (sub)gradient (a vector)!
- But only when we have *direct* access to the inner workings of f

$$\min_{x \in [a,b]} f(x), \quad \text{where } f \text{ is strictly quasiconvex}$$

Algorithm 1: Golden-section search

Input: $a < b, g = \frac{\sqrt{5}+1}{2}, \text{tol}$

```
1  $x_1 = a + (b - a)/g$ 
2  $x_2 = b - (b - a)/g$ 
3 while  $x_2 - x_1 > \text{tol}$  do
4   if  $f(x_2) > f(x_1)$  then
5      $b = x_2$ 
6      $x_2 = a + (b - a)/g$ 
7   else
8      $a = x_1$ 
9      $x_1 = b - (b - a)/g$ 
```

Fix the number of evaluations. Is there an “optimal” alg?

$$\inf_{\mathcal{A}} \sup_f \text{length of returned interval}$$

Key idea: recycle!

$$\min_{\lambda_2 \leq 1/2} \prod_{i=2}^N (1 - \lambda_i), \quad \text{s.t.} \quad \lambda_{n+1} = \frac{\lambda_n}{1 - \lambda_n} \wedge \frac{1 - 2\lambda_n}{1 - \lambda_n}$$

Solution: $\lambda_n = \frac{F_{n-1}}{F_{n+1}}$

Uniform Grid Search

Algorithm 2: Random pursuit

Input: \mathbf{w}_0 such that $\llbracket f \leq f(\mathbf{w}_0) \rrbracket$ is compact

```
1 for  $t = 1, 2, \dots$  do
2   choose normalized direction  $\mathbf{d}_t$  randomly
3    $\eta_t \leftarrow \operatorname{argmin}_{\eta \in \mathbb{R}} f(\mathbf{w}_t + \eta \mathbf{d}_t)$            // line search on chosen direction
4    $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta_t \mathbf{d}_t$ 
```

S. U. Stich, C. L. Müller, and B. Gärtner. "Optimization of Convex Functions with Random Pursuit". *SIAM Journal on Optimization*, vol. 23, no. 2 (2013), pp. 1284–1309, S. U. Stich, C. L. Müller, and B. Gärtner. "Variable metric random pursuit". *Mathematical Programming*, vol. 156 (2016), pp. 549–579.

If f is L_1 -smooth, then

$$\begin{aligned} f(\mathbf{w}_t + \eta_t \mathbf{d}_t) &\leq f(\mathbf{w}_t) + \eta_t \langle \mathbf{d}_t, \nabla f(\mathbf{w}_t) \rangle + \frac{L_1}{2} \eta_t^2 \\ &\leq f(\mathbf{w}_t) + \eta (\mathbf{w} - \mathbf{w}_t)^\top \mathbf{d}_t \mathbf{d}_t^\top \nabla f(\mathbf{w}_t) + \frac{L_1}{2} \eta^2 (\mathbf{w} - \mathbf{w}_t)^\top \mathbf{d}_t \mathbf{d}_t^\top (\mathbf{w} - \mathbf{w}_t)^\top \end{aligned}$$

- The above inequality is due to setting $\eta_t = \eta (\mathbf{w} - \mathbf{w}_t)^\top \mathbf{d}_t$ for some $\eta > 0$
- Using $\mathbb{E} \mathbf{d}_t \mathbf{d}_t^\top = \frac{1}{d} \mathbb{I}$ and assuming f is convex:

$$\begin{aligned} \mathbb{E} f(\mathbf{w}_t + \eta_t \mathbf{d}_t) &\leq f(\mathbf{w}_t) + \frac{\eta}{d} \langle \mathbf{w} - \mathbf{w}_t, \nabla f(\mathbf{w}_t) \rangle + \frac{\eta^2 L_1}{2d} \|\mathbf{w} - \mathbf{w}_t\|_2^2 \\ &\leq f(\mathbf{w}_t) + \frac{\eta}{d} [f(\mathbf{w}) - f(\mathbf{w}_t)] + \left(\frac{\eta}{d}\right)^2 \frac{d L_1}{2} \|\mathbf{w} - \mathbf{w}_t\|_2^2 \end{aligned}$$

- A simple induction (as in conditional gradient) yields:

$$\mathbb{E}[f(\mathbf{w}_t) - f(\mathbf{w})] \leq O\left(\frac{d L_1}{t+1}\right)$$

- A factor of dimension d worse

Finite Difference Approximation

$$\partial_j f(\mathbf{w}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{w} + t\mathbf{e}_j) - f(\mathbf{w})}{t}$$

- Choose small t often is enough, barring numerical cares
- Need to avoid doing this for every dimension
- Randomization may help!

Convolution

Definition: Convolution and Fourier transform

The convolution of two functions f and g is defined through integration:

$$(f * g)(\mathbf{w}) := \int_{\mathbf{z}} f(\mathbf{w} - \mathbf{z})g(\mathbf{z}) \, d\mathbf{z} = \int_{\mathbf{z}} f(\mathbf{z})g(\mathbf{w} - \mathbf{z}) \, d\mathbf{z} =: (g * f)(\mathbf{w}).$$

Recall the Fourier transform and its inverse:

$$(\mathcal{F}f)(\mathbf{w}^*) = \mathcal{F}f(\mathbf{w}^*) = \int_{\mathbf{w}} \exp(-2\pi i \langle \mathbf{w}, \mathbf{w}^* \rangle) f(\mathbf{w}) \, d\mathbf{w}$$

$$(\mathcal{F}^{-1}g)(\mathbf{w}) = \int_{\mathbf{w}^*} \exp(2\pi i \langle \mathbf{w}, \mathbf{w}^* \rangle) g(\mathbf{w}^*) \, d\mathbf{w}^*$$

$$\mathcal{F}(f * g) = \mathcal{F}f \cdot \mathcal{F}g, \quad \mathcal{F}\mathcal{F}^{-1} = \mathcal{F}^{-1}\mathcal{F} = \text{Id}, \quad \mathcal{F}f^{(\mathbf{k})} = (-2\pi i \mathbf{w}^*)^{\mathbf{k}} \mathcal{F}f$$

- Applying Fourier transform to the derivative of convolution:

$$\begin{aligned} \mathcal{F}(f * g)^{(\mathbf{k})} &= (-2\pi i \mathbf{w}^*)^{\mathbf{k}} \cdot \mathcal{F}(f * g) = [(-2\pi i \mathbf{w}^*)^{\mathbf{k}} \mathcal{F}f] \mathcal{F}g = \mathcal{F}(f^{(\mathbf{k})} * g) \\ &= \mathcal{F}(f * g^{(\mathbf{k})}) \end{aligned}$$

- Applying the inverse transform we obtain the formula of differentiating under the integral:

$$(f * g)^{(\mathbf{k})} = f^{(\mathbf{k})} * g = f * g^{(\mathbf{k})}$$

- This can in fact be the definition of the derivative (distribution) of f , using the derivative of some super smooth functions g !

Randomized Smoothing

Definition:

For a (vector-valued) function $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^c$ we define its **randomized smoothing** as

$$\mathbf{f}_\gamma(\mathbf{w}) = \mathbb{E}\mathbf{f}(\mathbf{w} + \gamma\boldsymbol{\epsilon}) = \mathbb{E}\mathbf{f}(\mathbf{w} - \gamma\boldsymbol{\epsilon}),$$

where $\boldsymbol{\epsilon}$ is some symmetric random noise with zero mean and identity covariance.

- Let p be the **probability density function** (pdf) of $\boldsymbol{\epsilon}$
- **Dilated density**: $p_\gamma(\mathbf{z}) = \frac{1}{\gamma^d} p(\frac{1}{\gamma}\mathbf{z})$
- We have point-wise convergence:

$$\mathbf{f}_\gamma = \mathbb{E}\mathbf{f}(\mathbf{w} - \gamma\boldsymbol{\epsilon}) = \mathbf{f} * p_\gamma, \text{ hence } \mathbf{f}_\gamma \rightarrow \mathbf{f} \text{ as } \gamma \rightarrow 0$$

- Intuitively expected, as the noise shrinks to 0, i.e. $p_\gamma \rightarrow \delta'_0$

Calculus for Randomized Smoothing

- The map $\mathbf{f} \mapsto \mathbf{f}_\gamma$ is linear
- If f is convex/concave, so is f_γ
- If f is convex, then $f_\gamma \geq f$
- If \mathbf{f} is L_0 -Lipschitz continuous (w.r.t. $\|\cdot\|_2$ say), so is \mathbf{f}_γ . Moreover,

$$\|\mathbf{f}_\gamma - \mathbf{f}\|_2 \leq \gamma L_0 \mathbb{E} \|\boldsymbol{\epsilon}\|_2 \leq \gamma L_0 \sqrt{\mathbb{E} \|\boldsymbol{\epsilon}\|_2^2} = \gamma L_0 \sqrt{d}$$

- If f is L_1 -smooth (w.r.t. $\|\cdot\|_2$ say), so is f_γ . Moreover,

$$f_\gamma - f \leq \frac{\gamma^2 L_1}{2} \mathbb{E} \|\boldsymbol{\epsilon}\|_2^2 = \frac{\gamma^2 L_1 d}{2},$$

whereas a two-sided bound holds if both $\pm f$ are L_1 -smooth.

Gradient approximation

- If $\pm f$ is L_1 -smooth, then $\|\nabla f_\gamma - \nabla f\|_o \leq \gamma L_1 \sqrt{d}$.
 - in fact, $\nabla f_\gamma = (\nabla f)_\gamma$, and $\|\nabla f\|_o \leq \|\nabla f_\gamma\|_o + \gamma L_1 \sqrt{d}$
- If $\pm f$ is L_2 -smooth, then $\|\nabla f_\gamma - \nabla f\|_o \leq \gamma^2 L_2 d/2$.
 - in fact, $\nabla f_\gamma = (\nabla f)_\gamma$ and $\nabla^2 f_\gamma = (\nabla^2 f)_\gamma$

Justifying the Name

Differentiating under the integral we obtain

$$f_{\gamma}^{(\mathbf{k})} := [f * p_{\gamma}]^{(\mathbf{k})} = f^{(\mathbf{k}-1)} * p_{\gamma}^{(1)}, \quad \nabla^k f_{\gamma}(\mathbf{w}) = \int \nabla^{k-1} f(\mathbf{w} - \mathbf{z}) \otimes \nabla p_{\gamma}(\mathbf{z}) \, d\mathbf{z}.$$

Therefore, if f is L_{k-1} -smooth, then f_{γ} is L_k -smooth, where

$$L_k \leq L_{k-1} \int \|\nabla p_{\gamma}(\mathbf{z})\|_2 \, d\mathbf{z} = \frac{L_{k-1}}{\gamma} \int \|\nabla p(\mathbf{z})\|_2 \, d\mathbf{z} = \frac{s L_{k-1}}{\gamma}$$

- $s := \mathbb{E} \|\nabla \ln p(\varepsilon)\|_2, \quad \varepsilon \sim p$
- f_{γ} is (at least) 1 degree more smoother than f

$$\begin{aligned}
\nabla f_\gamma(\mathbf{w}) &= \int f(\mathbf{w} - \mathbf{z}) \nabla p_\gamma(\mathbf{z}) \, d\mathbf{z} = \frac{1}{\gamma} \mathbb{E}[f(\mathbf{w} - \gamma \boldsymbol{\varepsilon}) \nabla \ln p(\boldsymbol{\varepsilon})] \\
&= -\frac{1}{\gamma} \mathbb{E}[f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) \nabla \ln p(\boldsymbol{\varepsilon})] \\
&= -\mathbb{E} \left[\frac{f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) - f(\mathbf{w})}{\gamma} \nabla \ln p(\boldsymbol{\varepsilon}) \right] \\
&= -\mathbb{E} \left[\frac{f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) - f(\mathbf{w} - \gamma \boldsymbol{\varepsilon})}{2\gamma} \nabla \ln p(\boldsymbol{\varepsilon}) \right]
\end{aligned}$$

When f is e.g. convex or an envelope function, we have the limit:

$$\begin{aligned}
\nabla f_0(\mathbf{w}) &:= -\mathbb{E}[f'(\mathbf{w}; \boldsymbol{\varepsilon}) \nabla \ln p(\boldsymbol{\varepsilon})], \quad \text{where} \quad f'(\mathbf{w}; \boldsymbol{\varepsilon}) := \lim_{\gamma \downarrow 0} [f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) - f(\mathbf{w})]/\gamma \\
&= -\mathbb{E}[\sigma_{\partial f(\mathbf{w})}(\boldsymbol{\varepsilon}) \nabla \ln p(\boldsymbol{\varepsilon})]
\end{aligned}$$

Needless to say, when f is actually differentiable, we have $\nabla f_0 = \nabla f$.

Gaussian Smoothing

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbb{I}), \quad i.e. \quad p(\boldsymbol{\varepsilon}) = (2\pi)^{d/2} \exp(-\|\boldsymbol{\varepsilon}\|_2^2/2)$$

- $-\nabla \ln p(\boldsymbol{\varepsilon}) = \boldsymbol{\varepsilon}$ and $s = \mathbb{E}\|\nabla \ln p(\boldsymbol{\varepsilon})\|_2 \leq \sqrt{d}$
- Conveniently, f_γ is in fact infinitely many times differentiable, e.g.

$$\nabla f_\gamma(\mathbf{w}) = \frac{1}{\gamma} \mathbb{E}[f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) \boldsymbol{\varepsilon}] = \mathbb{E} \left[\frac{f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) - f(\mathbf{w})}{\gamma} \boldsymbol{\varepsilon} \right] = \mathbb{E} \left[\frac{f(\mathbf{w} + \gamma \boldsymbol{\varepsilon}) - f(\mathbf{w} - \gamma \boldsymbol{\varepsilon})}{2\gamma} \boldsymbol{\varepsilon} \right]$$

- Requires f to be defined on entire \mathbb{R}^d

Uniform Smoothing

$$\boldsymbol{\varepsilon} \sim \text{Uniform}(K), \quad \text{i.e.} \quad p(\boldsymbol{\varepsilon}) = \begin{cases} 1/v_d, & \text{if } \boldsymbol{\varepsilon} \in K \\ 0, & \text{otherwise} \end{cases}$$

- v_d is the volume of the (symmetric, **isotropic**, i.e. $\mathbb{E}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top = \mathbb{I}$) compact set K
- Applying **Stokes' theorem**, $\nabla p(\boldsymbol{\varepsilon}) = \mathbf{1}_{\partial K} \cdot \mathbf{n}(\boldsymbol{\varepsilon})/v_d$, where $\mathbf{n}(\boldsymbol{\varepsilon})$ is the normal vector
- $s = u_{d-1}/v_d$ where u_{d-1} is the surface area of ∂K ; choose $\boldsymbol{\delta} \sim \text{Uniform}(\partial K)$:

$$\begin{aligned} \nabla f_\gamma(\mathbf{w}) &= -\frac{s}{\gamma} \mathbb{E}[f(\mathbf{w} + \gamma\boldsymbol{\delta})\mathbf{n}(\boldsymbol{\delta})] = -s \mathbb{E} \left[\frac{f(\mathbf{w} + \gamma\boldsymbol{\delta}) - f(\mathbf{w})}{\gamma} \mathbf{n}(\boldsymbol{\delta}) \right] \\ &= -s \mathbb{E} \left[\frac{f(\mathbf{w} + \gamma\boldsymbol{\delta}) - f(\mathbf{w} - \gamma\boldsymbol{\delta})}{2\gamma} \mathbf{n}(\boldsymbol{\delta}) \right] \end{aligned}$$

- Requires f to be defined (and bounded) over $C + \gamma K$.
- Let $K = \mathbf{B}_2(\mathbf{0}, \sqrt{d})$ we have $\mathbf{n}(\boldsymbol{\delta}) = -\sqrt{d}\boldsymbol{\delta}/\|\boldsymbol{\delta}\|_2$ and $s = \sqrt{d}$

Put Everything Together

- We optimize f_γ as a smoothed approximation of f
- We compute an **unbiased, stochastic** (sub)gradient of f_γ by
 1. $\hat{\partial}^1 f_\gamma(\mathbf{w}) = -\frac{1}{\gamma} f(\mathbf{w} + \gamma \epsilon) \cdot \nabla \ln p(\epsilon)$
 2. $\hat{\partial}^{1,0} f_\gamma(\mathbf{w}) = -\frac{f(\mathbf{w} + \gamma \epsilon) - f(\mathbf{w})}{\gamma} \cdot \nabla \ln p(\epsilon)$
 3. $\hat{\partial}^{1,1} f_\gamma(\mathbf{w}) = -\frac{f(\mathbf{w} + \gamma \epsilon) - f(\mathbf{w} - \gamma \epsilon)}{2\gamma} \cdot \nabla \ln p(\epsilon)$
 4. $\hat{\partial} f_0(\mathbf{w}) = -f'(\mathbf{w}; \epsilon) \cdot \nabla \ln p(\epsilon)$
- Except the last choice, only **require 1 or 2 evaluations of the function**
- Except the last choice, these stochastic (sub)gradients **in general are biased for f**
- We bound the second moment of the stochastic (sub)gradient
- We apply the stochastic GDA algorithm and obtain convergence towards f_γ
- We set γ appropriately so that we obtain convergence towards f

L_0 -Lipschitz Continuous and Convex

- If f is convex, then $f_\gamma \geq f$
- If f is L_0 -Lipschitz continuous (w.r.t. $\|\cdot\|_2$ say), so is f_γ . Moreover,

$$\|f_\gamma - f\|_2 \leq \gamma L_0 \mathbb{E} \|\epsilon\|_2 \leq \gamma L_0 \sqrt{\mathbb{E} \|\epsilon\|_2^2} = \gamma L_0 \sqrt{d}$$

- Thus, we obtain the approximation bound:

$$\mathbb{E}[f(\bar{\mathbf{w}}_t) - f(\mathbf{w})] - \gamma L_0 \sqrt{d} \leq \mathbb{E}[f_\gamma(\bar{\mathbf{w}}_t) - f_\gamma(\mathbf{w})]$$

- Using $\hat{\partial}^{1,0} f_\gamma$ we obtain

$$\mathbb{E}[f_\gamma(\bar{\mathbf{w}}_t) - f_\gamma(\mathbf{w})] \leq \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \eta_k^2 \cdot \mathbb{E} \|\hat{\partial}^{1,0} f_\gamma(\mathbf{w})\|_2^2}{2H_t}$$

- If f is L_0 -Lipschitz continuous, then using Gaussian smoothing:

$$\begin{aligned}\mathbb{E}\|\hat{\partial}^{1,0} f_\gamma(\mathbf{w})\|_2^2 &= \mathbb{E}\left\| -\frac{f(\mathbf{w}+\gamma\boldsymbol{\varepsilon})-f(\mathbf{w})}{\gamma} \cdot \nabla \ln p(\boldsymbol{\varepsilon}) \right\|_2^2 \\ &\leq L_0^2 \cdot \mathbb{E}\|\boldsymbol{\varepsilon}\|_2^4 \\ &\leq L_0^2 \cdot d(d+2) \leq L_0^2(d+1)^2\end{aligned}$$

- Setting $\gamma = \frac{\epsilon}{2L_0\sqrt{d}}$, $\eta_t = \frac{\text{diam}(C)}{(d+1)L_0\sqrt{t+1}}$ we have

$$\mathbb{E}[f(\bar{\mathbf{w}}_t) - f(\mathbf{w})] \leq \epsilon, \quad \text{if } t > \frac{4(d+1)^2}{\epsilon^2} [\text{diam}(C)L_0]^2,$$

which is d^2 times slower than running subgradient directly on f .

L_1 -smooth and convex

- If f is convex, then $f_\gamma \geq f$
- If f is L_1 -smooth (w.r.t. $\|\cdot\|_2$ say), so is f_γ . Moreover,

$$f_\gamma - f \leq \frac{\gamma^2 L_1}{2} \mathbb{E} \|\epsilon\|_2^2 = \frac{\gamma^2 L_1 d}{2}$$

- Thus, we obtain the approximation bound:

$$\mathbb{E}[f(\bar{\mathbf{w}}_t) - f(\mathbf{w})] - \frac{\gamma^2 L_1 d}{2} \leq \mathbb{E}[f_\gamma(\bar{\mathbf{w}}_t) - f_\gamma(\mathbf{w})]$$

- Using again $\hat{\partial}^{1,0} f_\gamma$ we obtain similarly

$$\mathbb{E}[f_\gamma(\bar{\mathbf{w}}_t) - f_\gamma(\mathbf{w})] \leq \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \eta_k^2 \cdot \mathbb{E} \|\hat{\partial}^{1,0} f_\gamma(\mathbf{w}_k)\|_2^2}{2H_t}$$

- If ∇f is L_1 -Lipschitz continuous:

$$\begin{aligned} \mathbb{E} \|\hat{\partial}^{1,0} f_\gamma(\mathbf{w})\|_2^2 &= \mathbb{E} \left\| -\frac{f(\mathbf{w} + \gamma \boldsymbol{\epsilon}) - f(\mathbf{w})}{\gamma} \cdot \nabla \ln p(\boldsymbol{\epsilon}) \right\|_2^2 \\ &\leq \mathbb{E} \left[\langle \nabla f(\mathbf{w}), \boldsymbol{\epsilon} \rangle + \frac{L_1 \gamma \|\boldsymbol{\epsilon}\|_2^2}{2} \right]^2 \|\boldsymbol{\epsilon}\|_2^2 \\ &\leq \frac{\gamma^2 L_1^2}{2} d(d+2)(d+4) + 2(d+2) \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

- With $\gamma = O\left(\frac{1}{d} \sqrt{\frac{\epsilon}{L_1}}\right)$ and $\eta_t \equiv O\left(\frac{1}{dL_1}\right)$, need $O\left(\frac{d}{\epsilon} L_1 \text{diam}^2(C)\right)$ many steps to obtain an ϵ -minimizer of f
- d times slower than running (projected) gradient directly on f

More Moment Bounds for Gaussian Smoothing

- If f is differentiable:

$$\begin{aligned}\mathbb{E}\|\hat{\partial}f_0(\mathbf{w})\|_2^2 &= \mathbb{E}\|\boldsymbol{\varepsilon}\|_2^4 \left\langle \frac{\boldsymbol{\varepsilon}}{\|\boldsymbol{\varepsilon}\|_2}, \nabla f(\mathbf{w}) \right\rangle^2 \\ &= \mathbb{E}\|\boldsymbol{\varepsilon}\|_2^4 \cdot \mathbb{E} \left\langle \frac{\boldsymbol{\varepsilon}}{\|\boldsymbol{\varepsilon}\|_2}, \nabla f(\mathbf{w}) \right\rangle^2 = (d+2)\|\nabla f(\mathbf{w})\|_2^2\end{aligned}$$

- If $\pm f$ is L_1^\pm -smooth:

$$\mathbb{E}\|\hat{\partial}^{1,1}f_\gamma(\mathbf{w})\|_2^2 \leq \frac{\gamma^2(L_1^+ + L_1^-)^2}{8}d(d+2)(d+4) + 2(d+2)\|\nabla f(\mathbf{w})\|_2^2$$

- If $\nabla^2 f$ is L_2 -Lipschitz continuous

$$\mathbb{E}\|\hat{\partial}^{1,1}f_\gamma(\mathbf{w})\|_2^2 \leq \frac{\gamma^4 L_2^2}{18}d(d+2)(d+4)(d+6) + 2(d+2)\|\nabla f(\mathbf{w})\|_2^2$$

