



Weierstrass Institute for
Applied Analysis and Stochastics



Random gradient free optimization: Bayesian view

Vladimir Spokoiny ,
WIAS, HU Berlin

April, 2020

1 Introduction

2 Bayesian optimization

3 Theoretical results for one-step posterior

- Conditions
- Examples
- Theoretical results

4 Details

- Properties of qMLE
- Posterior contraction
- Gaussian approximation of $v \mid Y$

Aim: An efficient procedure to minimize a convex function $f(\mathbf{x})$ without computing the gradient and Hessian.

[Nesterov and Spokoiny, 2017] offered a “gradient free” procedure which only relies on the directional derivative of $f(\mathbf{x})$

Procedure: with $\mathbf{u} \sim \mathcal{N}(0, B^{-1})$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k \frac{f(\mathbf{x}_k + \mu \mathbf{u}) - f(\mathbf{x})}{\mu} \cdot B \mathbf{u}.$$

The basic tool of study is the average

$$f_\mu(\mathbf{x}) = \mathbb{E} f(\mathbf{x} + \mu \mathbf{u}) = \frac{\int f(\mathbf{x} + \mu \mathbf{u}) e^{-\langle B \mathbf{u}, \mathbf{u} \rangle / 2} d\mathbf{u}}{\int e^{-\langle B \mathbf{u}, \mathbf{u} \rangle / 2} d\mathbf{u}}$$

Let \mathbf{Y} denote the observed random data, $\mathbf{Y} \sim \mathbb{P}$.

Model (DNN): $\mathbb{P} \in (\mathbb{P}_{\mathbf{v}}, \mathbf{v} \in \mathcal{Y} \subseteq \mathbb{R}^{\infty})$.

The log-likelihood function (negative fidelity)

$$L(\mathbf{v}) = L(\mathbf{Y}, \mathbf{v}) \stackrel{\text{def}}{=} \log \frac{d\mathbb{P}_{\mathbf{v}}}{d\mu_0}(\mathbf{Y}).$$

Training by MLE: maximizing the random function $L(\mathbf{v})$

$$\tilde{\mathbf{v}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} L(\mathbf{v}) = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \exp L(\mathbf{v}).$$

Target (the best parametric fit/ risk minimization):

$$\mathbf{v}^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E} L(\mathbf{v}) \neq \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E} \exp L(\mathbf{v}).$$

Also concavity of $L(\mathbf{v}) \neq$ concavity of $\exp L(\mathbf{v})$

- \mathcal{Y} is high or infinite dimensional
- $\nabla L(\mathbf{v})$ and $\nabla^2 L(\mathbf{v})$ hard to compute
- badly posed, need of regularization
- non-convex and non-smooth problem
- dimension reduction issue (drop-out)
- standard technique: SGD + backpropagation; efficient numerically but theoretical guarantees are hard to obtain

Bayesian (MCMC-type) methods with Gaussian priors yielding

- Efficient gradient and Hessian free procedure with second order accuracy

- Theoretical guarantees

1 Introduction

2 Bayesian optimization

3 Theoretical results for one-step posterior

- Conditions
- Examples
- Theoretical results

4 Details

- Properties of qMLE
- Posterior contraction
- Gaussian approximation of $v \mid Y$

In the **Bayes setup** \mathbf{v} is a random element, $\mathbf{v} \sim \Pi$ on the parameter set \mathcal{Y} , a **prior** with density $\Pi(\mathbf{v})$.

The **posterior** describes the conditional distribution of \mathbf{v} given \mathbf{Y}

$$\mathbf{v} | \mathbf{Y} \sim \frac{\exp\{L(\mathbf{v})\} \Pi(\mathbf{v})}{\int \exp\{L(\mathbf{v})\} \Pi(\mathbf{v}) d\mathbf{v}} = \frac{\exp\{L_{\Pi}(\mathbf{v})\}}{\int \exp\{L_{\Pi}(\mathbf{v})\} d\mathbf{v}}$$

with $L_{\Pi} = L(\mathbf{v}) + \log \Pi(\mathbf{v})$.

Prior Π induces a **penalty** $-\log \Pi(\mathbf{v})$ leading to **penalized MLE**

$$\tilde{\mathbf{v}}_{\Pi} = \underset{\mathbf{v}}{\operatorname{argmax}} L_{\Pi}(\mathbf{v}) = \underset{\mathbf{v}}{\operatorname{argmax}} \{L(\mathbf{v}) + \log \Pi(\mathbf{v})\}.$$

Formally, Bayes approach replaces the **max** of L_{Π} by the **soft-max**.

A Gaussian priors $\mathcal{N}(\bar{\mathbf{v}}, G^{-2})$ lead to quadratic penalization

$$\tilde{\mathbf{v}}_G = \operatorname{argmax}_{\mathbf{v}} L_G(\mathbf{v}) = \operatorname{argmax}_{\mathbf{v}} \{L(\mathbf{v}) - \|G(\mathbf{v} - \bar{\mathbf{v}})\|^2/2\};$$

cf. the Moreau–Yosida proximal-point method.

Posterior $\mathbf{v}_G \mid \mathbf{Y}$ is a random measure with the density

$$\mathbf{v}_G \mid \mathbf{Y} \sim \frac{\exp L_G(\mathbf{v})}{\int \exp L_G(\mathbf{v}) d\mathbf{v}} \propto \exp L_G(\mathbf{v}).$$

1. Select a starting prior $\Pi_0 = \mathcal{N}(\bar{\mathbf{v}}_0, G_0^{-2})$; Set $k = 0$;
2. Draw an independent sample $\mathbf{v}_1^{(k)}, \dots, \mathbf{v}_M^{(k)}$ from the prior $\mathcal{N}(\bar{\mathbf{v}}_k, G_k^{-2})$. For each $\mathbf{v}_m^{(k)}$, compute $L_k(\mathbf{v}_m^{(k)})$ with

$$L_k(\mathbf{v}) = L(\mathbf{v}) - \|G_k(\mathbf{v} - \bar{\mathbf{v}}_k)\|^2/2$$

and the corresponding weight

$$w_m^{(k)} = \exp L_k(\mathbf{v}_m^{(k)}).$$

3. Use the collection $(\mathbf{v}_m^{(k)}, w_m^{(k)})$ for $m \leq M$ (posterior) to build the next prior distribution $\Pi_{k+1} = \mathcal{N}(\bar{\mathbf{v}}_{k+1}, G_{k+1}^{-2})$.
4. Increase $k \rightarrow k + 1$ and repeat pp. 2 and 3 until convergence.

A prior is called **conjugated** if the corresponding posterior belongs to the same family of measures as prior. This reduces the step of computing the posterior to parameter update.

Our **main result** claims that the Gaussian prior for a regular parametric family yields a **nearly Gaussian posterior (nearly conjugated)**.

Therefore, it suffices to recompute the parameters of normal law.

It is natural and standard to use the **posterior mean**

$$\bar{\mathbf{v}}_{k+1} = \hat{\mathbf{v}}_k = \frac{1}{N_k} \sum_m w_m^{(k)} \mathbf{v}_m^{(k)}, \quad N_k = \sum_m w_m^{(k)}.$$

Alternatively, a robust (trimmed) mean could be used.

We suggest to apply $G_k^2 = \rho_k^{-2} G^2$ with a fixed G^2 .

The prior concentrates on the ρ_k -vicinity of $\bar{\mathbf{v}}_k$ and the same for posterior, so, ρ_k has the flavor of a step size.

Let $\mathbb{E}L_k(\mathbf{v})$ be (locally) concave. Define $D_k^2 = -\nabla^2 \mathbb{E}L_k(\bar{\mathbf{v}}_k)$ and consider a quadratic approximation of $L_k(\mathbf{v})$ around $\bar{\mathbf{v}}_k$.

As $\tilde{\mathbf{v}}_k = \operatorname{argmax}_{\mathbf{v}} L_k(\mathbf{v})$, it holds

$$L_k(\mathbf{v}) - L_k(\tilde{\mathbf{v}}_k) \approx -\|D_k(\mathbf{v} - \tilde{\mathbf{v}}_k)\|^2/2,$$

$$\nabla L_k(\tilde{\mathbf{v}}_k) = 0,$$

$$\nabla L_k(\bar{\mathbf{v}}_k) - \nabla L_k(\tilde{\mathbf{v}}_k) \approx -D_k^2(\bar{\mathbf{v}}_k - \tilde{\mathbf{v}}_k),$$

and hence $\bar{\mathbf{v}}_{k+1} = \tilde{\mathbf{v}}_k$ corresponds to the Newton update:

$$\bar{\mathbf{v}}_{k+1} - \bar{\mathbf{v}}_k \approx D_k^{-2} \nabla L_k(\bar{\mathbf{v}}_k).$$

The use of $\nabla L_k(\tilde{\mathbf{v}}_k) = 0$ yields with $\tilde{D}_k^2 = -\nabla^2 \mathbb{E} L_k(\tilde{\mathbf{v}}_k)$

$$\begin{aligned}
 & \frac{\int g(\mathbf{u}) \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u})\} d\mathbf{u}}{\int \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u})\} d\mathbf{u}} \\
 &= \frac{\int g(\mathbf{u}) \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G)\} d\mathbf{u}}{\int \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G)\} d\mathbf{u}} \\
 &= \frac{\int g(\mathbf{u}) \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}{\int \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}} \\
 &\approx \frac{\int g(\mathbf{u}) \exp\{-\|\tilde{D}_k \mathbf{u}\|^2/2\} d\mathbf{u}}{\int \exp\{-\|\tilde{D}_k \mathbf{u}\|^2/2\} d\mathbf{u}}
 \end{aligned}$$

hence $\mathbf{v}_k \mid \mathbf{Y} \rightsquigarrow \mathcal{N}(\tilde{\mathbf{v}}_k, \tilde{D}_k^{-2})$.

The procedure is gradient and Hessian free. Each step delivers a **posterior distribution** $\mathbf{v}_k \mid \mathbf{Y}$. We show that

$$\mathbf{v}_k \mid \mathbf{Y} \approx \mathcal{N}(\tilde{\mathbf{v}}_k, \tilde{D}_k^{-2}), \quad \tilde{D}_k^2 = D_k^2(\tilde{\mathbf{v}}_k)$$

and, in particular,

- $\mathbf{v}_k \mid \mathbf{Y}$ concentrates on an elliptic vicinity of $\tilde{\mathbf{v}}_k$
- posterior mean $\hat{\mathbf{v}}_k$ is a proxi for $\tilde{\mathbf{v}}_k$ and can be computed from Bayesian sampling: $\hat{\mathbf{v}}_k = N_k^{-1} \sum_m \mathbf{v}_m^{(k)} w_m^{(k)}$;
- Information from all previous steps can be incorporated in the Gaussian prior $\mathcal{N}(\bar{\mathbf{v}}_k, G_k^{-2})$;
- Prior precision matrix G_k^2 can be used to manipulate with the step size and effective dimension.

Even in the parametric case, the value \mathbf{v}^* can be estimated with accuracy $n^{-1/2}$ at best.

Therefore, no sense to continue the procedure if $D_k^{-2} \ll n^{-1}$.

If ρ_k decreases exponentially, only $\log n$ steps are required.

1 Introduction

2 Bayesian optimization

3 Theoretical results for one-step posterior

- Conditions
- Examples
- Theoretical results

4 Details

- Properties of qMLE
- Posterior contraction
- Gaussian approximation of $v \mid Y$

(E) The *stochastic component* $\zeta(\mathbf{v}) = L(\mathbf{v}) - \mathbb{E}L(\mathbf{v})$ of the process $L(\mathbf{v})$ is *linear* in \mathbf{v} :

$$\nabla \zeta \equiv \nabla \zeta(\mathbf{v}).$$

(ED₀) There exist a positive symmetric matrix V , and constants $g > 0$, $\nu_0 \geq 1$ such that $\text{Var}(\nabla \zeta) \leq V^2$ and

$$\sup_{\mathbf{u} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\langle \mathbf{u}, \nabla \zeta \rangle}{\|V\mathbf{u}\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g.$$

- (\mathcal{L}) The set \mathcal{Y} is open and convex in \mathbb{R}^p . For each k , the function $\mathbb{E}L_k(\mathbf{v})$ is **concave** in $\mathbf{v} \in \mathcal{Y}$.
- (\mathcal{L}_0) Define for each $\mathbf{v} \in \mathcal{Y}^\circ$, and any $\mathbf{u} \in \mathbb{R}^p$, the directional derivative

$$\delta_m(\mathbf{v}, \mathbf{u}) \stackrel{\text{def}}{=} \frac{1}{m!} \frac{d^m}{dt^m} \mathbb{E}L(\mathbf{v} + t\mathbf{u}) \Big|_{t=0}, \quad m = 3, 4.$$

The functions $\delta_3(\mathbf{v}, \mathbf{u})$ and $\delta_4(\mathbf{v}, \mathbf{u})$ are well defined and with $D^2(\mathbf{v}) = -\nabla^2 \mathbb{E}L(\mathbf{v})$

$$\omega_m(\mathbf{v}) = \sup_{\mathbf{u}: \|D(\mathbf{v})\mathbf{u}\| \leq r} \frac{\delta_m(\mathbf{v}, \mathbf{u})}{\|D(\mathbf{v})\mathbf{u}\|^2} \leq \frac{1}{3}.$$

Define

$$\mathbb{F}(\mathbf{v}) = -\nabla^2 \mathbb{E}L(\mathbf{v}), \quad \mathbb{F}_G(\mathbf{v}) = \mathbb{F}(\mathbf{v}) + G^2.$$

(V|G) Signal-noise:

$$B_{V|G}(\mathbf{v}) \stackrel{\text{def}}{=} \mathbb{F}_G^{-1/2}(\mathbf{v}) V^2 \mathbb{F}_G^{-1/2}(\mathbf{v})$$

with V^2 from (ED_0) . There are fixed constants $\lambda_{V|G}$ and $\mathfrak{p}_{V|G}$ such that

$$\text{tr } B_{V|G}(\mathbf{v}) \leq \mathfrak{p}_{V|G}, \quad \|B_{V|G}\| \leq \lambda_{V|G}, \quad \mathbf{v} \in \mathcal{I}^\circ.$$

(D|G) The **effective dimension**: for a fixed constant C

$$\mathfrak{p}_G(\mathbf{v}) \stackrel{\text{def}}{=} \text{tr} \{ \mathbb{F}(\mathbf{v}) \mathbb{F}_G^{-1}(\mathbf{v}) \} \leq C, \quad \mathbf{v} \in \mathcal{I}^\circ.$$

Consider the model

$$\mathbf{Y} = A(\mathbf{f}) + \sigma \boldsymbol{\varepsilon} \in \mathcal{Y}^d$$

with a known non-linear operator $A: \mathcal{X} \rightarrow \mathcal{Y}^d$ for Hilbert spaces \mathcal{X}, \mathcal{Y} and a discretized subspace $\mathcal{Y}^d \subset \mathcal{Y}$.

Examples include

- PDE with elliptic operators, [Nickl et al., 2018]
- Schrödinger equation, [Nickl, 2017]
- Calderón equation, [Abraham and Nickl, 2019]

$$\mathbf{f} \mid \mathbf{Y} \propto \exp L_G(\mathbf{f}),$$
$$L_G(\mathbf{f}) = -\frac{1}{2\sigma^2} \|\mathbf{Y} - A(\mathbf{f})\|^2 - \frac{1}{2} \|G\mathbf{f}\|^2.$$

Let observations Y_i, X_i follow the model

$$Y_i \sim P_{f(X_i)} \in (P_u),$$

where (P_u) is an exponential family with a log-density

$$\ell(y, u) = \log p(y, u) = C(u)y - B(u);$$

$C(u)$ is an increasing and $B(u)$ is a convex function.

Include binary-response, Poissonian regression, Cox regression, reliability and extreme values, ...

PA (DNN): $f(x) = f(x, \mathbf{v}^*)$. Yields the log-likelihood

$$L(\mathbf{v}) = \sum_{i=1}^n \ell(Y_i, f(X_i, \mathbf{v})) = \sum_{i=1}^n Y_i \{C(f(X_i, \mathbf{v})) - B(f(X_i, \mathbf{v}))\}.$$

■ MAP and posterior mean

$$\tilde{\mathbf{v}}_{\text{MAP}} = \underset{\mathbf{v}}{\operatorname{argmax}} \exp L_G(\mathbf{v}) = \tilde{\mathbf{v}}_G$$

$$\bar{\mathbf{v}}_G \stackrel{\text{def}}{=} \frac{\int \mathbf{v} \exp L_G(\mathbf{v}) d\mathbf{v}}{\int \exp L_G(\mathbf{v}) d\mathbf{v}};$$

■ Concentration set \mathcal{A} :

$$\rho(\mathcal{A}) \stackrel{\text{def}}{=} \frac{\int_{\mathcal{A}^c} \exp L_G(\mathbf{v}) d\mathbf{v}}{\int_{\mathcal{A}} \exp L_G(\mathbf{v}) d\mathbf{v}}$$

■ Credible sets $\mathcal{A}(\alpha)$

$$\mathbb{P}(\mathbf{v}_G \in \mathcal{A}(\alpha) \mid \mathbf{Y}) = 1 - \alpha;$$

■ Elliptic credible sets $\mathcal{A}_{Q|G}(\alpha) = \{\mathbf{v} : \|Q(\mathbf{v}_G - \tilde{\mathbf{v}}_G)\| \leq z_\alpha\}$.

Penalized MLE (pMLE): with $L_G(\mathbf{v}) = L(\mathbf{v}) - \|G\mathbf{v}\|^2/2$

$$\tilde{\mathbf{v}}_G = \underset{\mathbf{v}}{\operatorname{argmax}} L_G(\mathbf{v}), \quad \mathbf{v}_G^* = \underset{\mathbf{v}}{\operatorname{argmax}} \mathbb{E}L_G(\mathbf{v})$$

Full information matrix (operator)

$$\mathbf{F}(\mathbf{v}) = -\nabla^2 \mathbb{E}L(\mathbf{v}),$$

$$\mathbf{F}_G(\mathbf{v}) = -\nabla^2 \mathbb{E}L_G(\mathbf{v}) = \mathbf{F}(\mathbf{v}) + G^2,$$

and

$$D_G^2 = \mathbf{F}_G(\mathbf{v}_G^*), \quad \tilde{D}_G^2 = \mathbf{F}_G(\tilde{\mathbf{v}}_G).$$

Effective dimension

$$p_G(\mathbf{v}) = \operatorname{tr}\left(\mathbf{F}_G^{-1}(\mathbf{v})\mathbf{F}(\mathbf{v})\right), \quad p_G = p_G(\mathbf{v}_G^*), \quad \tilde{p}_G = p_G(\tilde{\mathbf{v}}_G).$$

Fisher expansion for pMLE: on a set Ω of high probability

$$\|D_G(\tilde{v}_G - \mathbf{v}_G^*) - D_G^{-1}\nabla\zeta\|^2 \lesssim \sqrt{\frac{p_G}{n}} \|D_G^{-1}\nabla\zeta\|^2$$

with $\nabla\zeta = \nabla L_G(\mathbf{v}_G^*)$.

Wilks expansion

$$\left| 2L_G(\tilde{v}_G) - 2L_G(\mathbf{v}_G^*) - \|D_G^{-1}\nabla\zeta\|^2 \right| \lesssim \sqrt{\frac{p_G}{n}} \|D_G^{-1}\nabla\zeta\|^2,$$

$$\left| 2L_G(\tilde{v}_G) - 2L_G(\mathbf{v}) - \|D_G(\tilde{v}_G - \mathbf{v})\|^2 \right| \lesssim \sqrt{\frac{p_G}{n}} \|D_G^{-1}\nabla\zeta\|^2.$$

$$\sup_{A \in \mathcal{B}_s(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| \lesssim \frac{\mathfrak{p}_G^3}{n}$$

$$\sup_{A \in \mathcal{B}(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| \lesssim \sqrt{\frac{\mathfrak{p}_G^3}{n}}$$

where $\mathcal{B}(\mathbb{R}^p)$ stands for all Borel sets while $\mathcal{B}_s(\mathbb{R}^p)$ all **centrally symmetric** Borel sets in \mathbb{R}^p .

- rate of estimation of pMLE
- posterior concentration and contraction rate
- use of posterior mean in place of MAP
- credible sets as frequentist confidence sets
- prior impact
- empirical or full Bayes approach for prior selection

All for

- finite samples
- explicit error terms via effective dimension instead of full parameter dimension

Most of results requires the upper bound on the **effective dimension**

$$p_G = \text{tr}(\mathbb{F}_G^{-1} \mathbb{F})$$

$$p_G \ll n$$

However, the main result on Gaussian approximation of the posterior only valid under

$$p_G^3 \ll n$$

1 Introduction

2 Bayesian optimization

3 Theoretical results for one-step posterior

- Conditions
- Examples
- Theoretical results

4 Details

- Properties of qMLE
- Posterior contraction
- Gaussian approximation of $v \mid Y$

By (E) , the stochastic component $\zeta(\mathbf{v}) = L(\mathbf{v}) - \mathbb{E}L(\mathbf{v})$ is linear in \mathbf{v} and $\nabla\zeta = \nabla\zeta(\mathbf{v})$ does not depend on \mathbf{v} .

Theorem Under condition (ED_0) , there exists a random set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - Ce^{-x}$ such that on this set

$$\|D_G^{-1}\nabla\zeta\| \leq z(B_{V|G}, \mathbf{x}),$$

where $B_{V|G} = D_G^{-1}V^2D_G^{-1}$ and

$$z(B_{V|G}, \mathbf{x}) = \sqrt{\text{tr } B_{V|G}} + \sqrt{2\mathbf{x}\lambda_{\max}(B_{V|G})}.$$

Let $\mathbf{v}_G^* = \operatorname{arginf}_{\mathbf{v}} \mathbb{E} L_G(\mathbf{v})$ and $D_G^2 = D^2(\mathbf{v}_G^*) + G^2$.

Theorem Let $\|D_G^{-1} \nabla \zeta\| \leq z(B_{V|G}, \mathbf{x})$ on a random set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - e^{-x}$. Let

$$\mathcal{A}_G(\mathbf{r}_G) \stackrel{\text{def}}{=} \{\mathbf{v} : \|D_G(\mathbf{v} - \mathbf{v}_G^*)\| \leq \mathbf{r}_G\}$$
$$(1 - \rho)\mathbf{r}_G \geq z(B_{V|G}, \mathbf{x}).$$

Then on $\Omega(\mathbf{x})$

$$\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}_G^*)\| \leq \mathbf{r}_G.$$

Local set

$$\mathcal{A}_G(\mathbf{r}_G) \stackrel{\text{def}}{=} \{\mathbf{v}: \|D_G(\mathbf{v} - \mathbf{v}_G^*)\| \leq \mathbf{r}_G\}.$$

Use **local smoothness** of $-\mathbb{E}L_G(\mathbf{v})$ to show that for $\mathbf{v} \in \mathcal{A}_G(\mathbf{r}_G)$

$$\begin{aligned} -\{\mathbb{E}L_G(\mathbf{v}) - \mathbb{E}L_G(\mathbf{v}_G^*)\} &\approx \|D_G(\mathbf{v} - \mathbf{v}_G^*)\|^2/2, \\ -\nabla\mathbb{E}L_G(\mathbf{v}) &\approx D_G^2(\mathbf{v} - \mathbf{v}_G^*) \end{aligned}$$

Use **convexity** of $\mathbb{E}L_G(\mathbf{v})$ to show that

$$\|\nabla\mathbb{E}L_G(\mathbf{v})\| \geq \mathbf{r}_G, \quad \mathbf{v} \notin \mathcal{A}_G(\mathbf{r}_G).$$

Use $\|D_G^{-1}\nabla\zeta\| \leq z(B_{V|G}, \mathbf{x})$ to show

$$\nabla L(\mathbf{v}) = \nabla\mathbb{E}L(\mathbf{v}) + \nabla\zeta \neq 0, \quad \mathbf{v} \notin \mathcal{A}_G(\mathbf{r}_G).$$

Concentration sets of the posterior in nonparametric models using

- Empirical process theory, large deviations of the log-likelihood and covering numbers and chaining arguments
- small ball probability
- local smoothness

See e.g.

- Ghoshal, S., Ghosh, J. K., van der Vaart, A. W. (2000) Convergence rates of posterior distributions. *Ann.Statist.*, 28, 500–531.
- van der Vaart, A. W., van Zanten, J. H. (2008) Rates of contraction of posterior distributions based on Gaussian process priors. *Ann. Statist.*, 36, 1031–1508.

Define the elliptic set $\{\mathbf{u}: \|\tilde{D}\mathbf{u}\| \leq r_0\}$ with $\tilde{D}^2 = D^2(\tilde{\mathbf{v}}_G)$.
Consider the random quantity

$$\rho(r_0) \stackrel{\text{def}}{=} \frac{\int_{\|\tilde{D}\mathbf{u}\| > r_0} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u})\} d\mathbf{u}}{\int_{\|\tilde{D}\mathbf{u}\| \leq r_0} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u})\} d\mathbf{u}}.$$

Theorem Let, for some fixed values r_0 and $x > 0$, it hold

$$C_0 r_0 \geq 2\sqrt{p_G(\mathbf{v})} + \sqrt{x}, \quad \mathbf{v} \in \mathcal{Y}^\circ.$$

Then, on the random set $\Omega(x)$ from Theorem 31, with $\tilde{p}_G = p_G(\tilde{\mathbf{v}}_G)$

$$\rho(r_0) \leq \exp\{-(\tilde{p}_G + x)/2\}.$$

Let $\tilde{\mathbf{v}}_G = \operatorname{arginf}_{\mathbf{v}} L_G(\mathbf{v})$. The use of $\nabla L_G(\tilde{\mathbf{v}}_G) = 0$ allows to represent

$$\begin{aligned}\rho(\mathbf{r}_0) &= \frac{\int_{\|\tilde{D}\mathbf{u}\| > \mathbf{r}_0} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G)\} d\mathbf{u}}{\int_{\|\tilde{D}\mathbf{u}\| \leq \mathbf{r}_0} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G)\} d\mathbf{u}} \\ &= \frac{\int_{\|\tilde{D}\mathbf{u}\| > \mathbf{r}_0} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}{\int_{\|\tilde{D}\mathbf{u}\| \leq \mathbf{r}_0} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}\end{aligned}$$

Fix $\mathbf{v} \in \mathcal{A}_G(\mathbf{r}_G)$. Consider $f(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$. As $\zeta(\mathbf{v}) = L(\mathbf{v}) - \mathbb{E}L(\mathbf{v})$ is linear in \mathbf{v} , it holds

$$\begin{aligned} L_G(\mathbf{v} + \mathbf{u}) - L_G(\mathbf{v}) - \langle \nabla L_G(\mathbf{v}), \mathbf{u} \rangle \\ = f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle. \end{aligned}$$

Therefore, it suffices to bound the ratio

$$\rho(\mathbf{v}) \stackrel{\text{def}}{=} \frac{\int_{\mathcal{U}^c} \exp\{f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle\} d\mathbf{u}}{\int_{\mathcal{U}} \exp\{f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle\} d\mathbf{u}}$$

for the elliptic set $\mathcal{U} = \mathcal{U}(\mathbf{v}, \mathbf{r}_0) = \{\mathbf{u}: \|D(\mathbf{v})\mathbf{u}\| \leq \mathbf{r}_0\}$ uniformly in \mathbf{v} from the set $\{\mathbf{v}: \|D_G(\mathbf{v} - \mathbf{v}_G^*)\| \leq \mathbf{r}_G\}$; see Theorem 31.

First we present some bounds for the denominator of $\rho(\mathbf{v})$. Local smoothness of $f(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$ implies

$$\begin{aligned} & \int_{\mathcal{U}} \exp\{f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle\} d\mathbf{u} \\ & \approx \int_{\mathcal{U}} \exp\left(-\frac{\|D_G(\mathbf{v})\mathbf{u}\|^2}{2}\right) d\mathbf{u}, \end{aligned}$$

Moreover,

$$\frac{\det D_G(\mathbf{v})}{(2\pi)^{p/2}} \int_{\mathcal{U}} \exp\left(-\frac{\|D_G(\mathbf{v})\mathbf{u}\|^2}{2}\right) d\mathbf{u} = \mathbb{P}(\|D(\mathbf{v})D_G^{-1}(\mathbf{v})\boldsymbol{\gamma}\| \leq r_0)$$

for $\boldsymbol{\gamma} \sim \mathcal{N}(0, I_p)$. The choice $r_0 \geq \sqrt{p_G(\mathbf{v})} + \sqrt{2x}$ yields

$$\mathbb{P}(\|D(\mathbf{v})D_G^{-1}(\mathbf{v})\boldsymbol{\gamma}\| \leq r_0) \geq 1 - e^{-x}.$$

Step 3: $\int_{\mathcal{U}_c} \exp\{f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle\} d\mathbf{u}$

$f(\mathbf{v}) = \mathbb{E}L(\mathbf{v})$ is concave and $-\langle \nabla^2 f(\mathbf{v})\mathbf{u}, \mathbf{u} \rangle = \|D(\mathbf{v})\mathbf{u}\|^2$.

For any \mathbf{u} with $\|D(\mathbf{v})\mathbf{u}\| = r > r_0$

$$\begin{aligned} f(\mathbf{v} + \mathbf{u}) - f(\mathbf{v}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle - \|G\mathbf{u}\|^2/2 \\ \leq -C_0(\|D(\mathbf{v})\mathbf{u}\|r_0 - r_0^2/2) - \|G\mathbf{u}\|^2/2 \\ = -C_0(\|D(\mathbf{v})\mathbf{u}\|r_0 - r_0^2/2) - \|D_G(\mathbf{v})\mathbf{u}\|^2/2 + \|D(\mathbf{v})\mathbf{u}\|^2/2. \end{aligned}$$

with $C_0 \geq 1/2$ and $D_G^2(\mathbf{v}) = D^2(\mathbf{v}) + G^2$.

Now we can use the result about Gaussian integrals:

$$\begin{aligned} \frac{\det D_G}{(2\pi)^{p/2}} \int_{\|D\mathbf{u}\| \geq r_0} \exp\left\{-\left(\|D\mathbf{u}\|r_0 - r_0^2/2 - \|D_G\mathbf{u}\|^2 + \|D\mathbf{u}\|^2\right)/2\right\} d\mathbf{u} \\ \leq C e^{-(p_G(\mathbf{v})+x)/2}. \end{aligned}$$

- van der Vaart, A. W. , van Zanten, J. H. (2008)
- Leahu, H. (2011). Gaussian models
- Castillo, I. , Nickl, R. (2013, 2014). General results.
- Panov and Sp (2015). Semiparametric problems

Contraction results and local quadratic approximation of the log-likelihood. The Gaussian approximation in a weak sense.

Theorem Suppose that It holds on $\Omega(\mathbf{x})$ with $\tilde{D}_G^2 = D_G^2(\tilde{\mathbf{v}}_G)$

$$\sup_{A \in \mathcal{B}_s(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| \leq \mathbf{C} \diamond(\mathbf{r}_0)$$

$$\sup_{A \in \mathcal{B}(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| \leq \mathbf{C} \delta_3(\mathbf{r}_0)$$

with

$$\delta_3(\mathbf{r}_0) \lesssim \frac{\mathbf{r}_0^3}{\sqrt{n}} \lesssim \frac{\mathbf{p}_G^{3/2}}{\sqrt{n}},$$

$$\diamond(\mathbf{r}_0) \lesssim \frac{\mathbf{r}_0^6}{n} \lesssim \frac{\mathbf{p}_G^3}{n}.$$







Fix any centrally symmetric set A . First we restrict the posterior probability to the set $\tilde{\mathcal{A}}(\mathbf{r}_0) = \{\mathbf{u} : \|\tilde{D}\mathbf{u}\| \leq \mathbf{r}_0\}$. Then we apply the quadratic approximation of the log-likelihood function $L(\mathbf{v})$. Denote $A(\mathbf{r}_0) = A \cap \tilde{\mathcal{A}}(\mathbf{r}_0)$. Obviously, $A(\mathbf{r}_0)$ is centrally symmetric as well. Further,

$$\begin{aligned} \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) &= \frac{\int_A \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u})\} d\mathbf{u}}{\int_{\mathbb{R}^p} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u})\} d\mathbf{u}} \\ &\leq \frac{\int_{A(\mathbf{r}_0)} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}{\int_{\tilde{\mathcal{A}}(\mathbf{r}_0)} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}} + \rho(\mathbf{r}_0). \end{aligned}$$

Fix $\mathbf{v} \in \mathcal{A}_G(\mathbf{r}_G)$. Then for $A \subset \mathcal{U}$ centrally symmetric

$$\begin{aligned} & \int_A \exp\{f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle\} d\mathbf{u} \\ & \geq (1 - \diamond(\mathbf{r}_0)) \int_A \exp\left(-\frac{\|D_G(\mathbf{v})\mathbf{u}\|^2}{2}\right) d\mathbf{u}, \\ & \int_A \exp\{f(\mathbf{v} + \mathbf{u}) - f(\mathbf{u}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle\} d\mathbf{u} \\ & \leq (1 + \diamond(\mathbf{r}_0)) \int_A \exp\left(-\frac{\|D_G(\mathbf{v})\mathbf{u}\|^2}{2}\right) d\mathbf{u}, \end{aligned}$$

where $\diamond(\mathbf{r}_0) = 4\delta_3^2 + 4\delta_4$.

-  Abraham, K. and Nickl, R. (2019).
On statistical Calderón problems.
-  Asi, H. and Duchi, J. C. (2019).
Stochastic (Approximate) Proximal Point Methods: Convergence, Optimality, and Adaptivity.
SIAM Journal on Optimization, 29(3):2257–2290.
-  Nesterov, Y. and Spokoiny, V. (2017).
Random gradient-free minimization of convex functions.
Foundations of Computational Mathematics, 17(2):527–566.
-  Nickl, R. (2017).
Bernstein - von Mises theorems for statistical inverse problems I: Schrödinger equation.
arXiv e-prints, page arXiv:1707.01764.
-  Nickl, R. and Söhl, J. (2019).
Bernstein – von Mises theorems for statistical inverse problems II: compound Poisson processes.
Electronic Journal of Statistics, 13(2):3513–3571.
-  Nickl, R., van de Geer, S. A., and Wang, S. (2018).
Convergence rates for Penalised Least Squares Estimators in PDE-constrained regression problems.



Spokoiny, V. and Panov, M. (2019).

Accuracy of Gaussian approximation in nonparametric Bernstein – von Mises Theorem.

preprint arXiv:https://arxiv.org/abs/1910.06028.



Trabs, M. (2018).

Bayesian inverse problems with unknown operators.

Inverse Problems, 34(8):085001.