



Weierstraß-Institut für
Angewandte Analysis und Stochastik



Gaussian Variational Inference

Vladimir Spokoiny ,
WIAS, HU Berlin

10. Oktober 2024

1 Gaussian Variational Inference

2 A basic lemma

- Fourth order approximation
- Quadratic penalization

3 Solution to VI problem

4 Optimization vs sampling

Let $\mathbb{P}_f \sim \exp f(\mathbf{x})$. Denote by $\mathbb{N}_{\mathbf{x}, \mathbb{Z}}$ the Gaussian measure with the mean \mathbf{x} and covariance \mathbb{Z}^{-1} , i.e. $\mathbb{N}_{\mathbf{x}, \mathbb{Z}} \stackrel{\text{def}}{=} \mathcal{N}(\mathbf{x}, \mathbb{Z}^{-1})$.

$$\text{Gauss VI: } (\mathbf{x}_{\text{VI}}, \mathbb{Z}_{\text{VI}}) = \underset{\mathbf{x}, \mathbb{Z}}{\operatorname{arginf}} \mathcal{K}(\mathbb{N}_{\mathbf{x}, \mathbb{Z}} \| \mathbb{P}_f).$$

Natural candidates:

1. **Laplace:** $\mathbf{x}_{\text{VI}} \approx \operatorname{argmax} f(\mathbf{x})$, $\mathbb{Z}_{\text{VI}} \approx -\nabla^2 f(\mathbf{x}^*)$;
2. **Moments:** $\mathbf{x}_{\text{VI}} \approx \mathbb{E}_f \mathbf{X}$, $\mathbb{Z}_{\text{VI}}^{-1} \approx \operatorname{Var}_f(\mathbf{X})$.

Let $\mathbb{P}_f \sim \exp f(\mathbf{x})$. Denote by $\mathbb{N}_{\mathbf{x}, \mathbb{Z}}$ the Gaussian measure with the mean \mathbf{x} and covariance \mathbb{Z}^{-1} , i.e. $\mathbb{N}_{\mathbf{x}, \mathbb{Z}} \stackrel{\text{def}}{=} \mathcal{N}(\mathbf{x}, \mathbb{Z}^{-1})$.

$$\text{Gauss VI: } (\mathbf{x}_{\text{VI}}, \mathbb{Z}_{\text{VI}}) = \underset{\mathbf{x}, \mathbb{Z}}{\operatorname{arginf}} \mathcal{H}(\mathbb{N}_{\mathbf{x}, \mathbb{Z}} \| \mathbb{P}_f).$$

Natural candidates:

1. **Laplace:** $\mathbf{x}_{\text{VI}} \approx \operatorname{argmax} f(\mathbf{x})$, $\mathbb{Z}_{\text{VI}} \approx -\nabla^2 f(\mathbf{x}^*)$;
2. **Moments:** $\mathbf{x}_{\text{VI}} \approx \mathbb{E}_f \mathbf{X}$, $\mathbb{Z}_{\text{VI}}^{-1} \approx \operatorname{Var}_f(\mathbf{X})$.

[Katsevich and Rigollet, 2023] argued for (2).

- [David M. Blei and McAuliffe, 2017] Variational Inference: A review for statisticians
- [Zhang and Gao, 2020] Convergence rates of variational posterior distributions
- [Wang and Blei, 2019] Frequentist consistency of variational Bayes
- [Han and Yang, 2019] Statistical inference in mean-field variational Bayes
- [Challis and Barber, 2013] Gaussian Kullback-Leibler approximate inference
- [Alquier and Ridgway, 2020] Concentration of tempered posteriors and of their variational approximations
- [Lambert et al., 2023] Variational inference via Wasserstein gradient flows

The VI approach assumes minimizing of the KL-divergence $\mathcal{K}(\mathbf{N}_{\mathbf{x}, \mathbb{Z}} \parallel \mathbf{P}_f)$ over all feasible \mathbf{x}, \mathbb{Z} . Here we rewrite this problem in terms of local parameters \mathbf{a} and S .

Lemma

For any \mathbf{x} and any \mathbb{Z} , it holds

$$\mathcal{K}(\mathbf{P}_{\mathbf{x}, \mathbb{Z}} \parallel \mathbf{P}_f) = \mathbf{C} + \frac{1}{2} \log \det(\mathbb{Z}^{-1}) - \frac{p}{2} - \mathbb{E} f(\mathbf{x} + \boldsymbol{\gamma}_{\mathbb{Z}}).$$

with \mathbf{C} depending on f and p only.

With $C_f \stackrel{\text{def}}{=} \log \int e^{f(\bar{x}+u)} d\mathbf{u}$ and $C_p = (2\pi)^{-p/2}$, for any $\mathbf{u} \in \mathbb{R}^p$

$$\frac{d\mathbb{P}_f}{d\mathbf{u}}(\mathbf{x} + \mathbf{u}) = e^{-C_f} e^{f(\mathbf{x}+\mathbf{u})},$$

$$\frac{d\mathbb{P}_{\mathbf{x}, \mathbb{Z}}}{d\mathbf{u}}(\mathbf{x} + \mathbf{u}) = C_p \det(\mathbb{Z}^{1/2}) e^{-\|\mathbb{Z}^{1/2}\mathbf{u}\|^2/2}.$$

This yields with $\gamma_{\mathbb{Z}} \sim \mathcal{N}(0, \mathbb{Z}^{-1})$ and $\gamma \sim \mathcal{N}(0, \mathbb{I}_p)$

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, \mathbb{Z}} \log \frac{d\mathbb{P}_{\mathbf{x}, \mathbb{Z}}}{d\mathbb{P}_f} \\ = C_f + \log C_p - \mathbb{E} f(\mathbf{x} + \gamma_{\mathbb{Z}}) - \frac{1}{2} \mathbb{E} \|\gamma\|^2 - \frac{1}{2} \log \det(\mathbb{Z}^{-1}), \end{aligned}$$

and the result follows in view of $\mathbb{E} \|\gamma\|^2 = p$.

With $\mathbb{F} = -\nabla^2 f(\bar{\mathbf{x}})$, represent \mathbb{Z} in the form

$$\mathbb{Z}^{-1/2} = \mathbb{F}^{-1/2} + S \quad \text{or} \quad \mathbb{F}^{1/4} \mathbb{Z}^{-1/2} \mathbb{F}^{1/4} = \mathbb{I}_p + \mathbb{F}^{1/4} S \mathbb{F}^{1/4}.$$

A vicinity of \mathbb{F} using Kullback-Leibler divergence $\mathcal{K}(\mathbf{N}_{\bar{\mathbf{x}},\mathbb{F}} \parallel \mathbf{N}_{\bar{\mathbf{x}},\mathbb{Z}})$.

Lemma

Let $\mathbb{Z}^{-1/2} = \mathbb{F}^{-1/2} + S$ and $U = \mathbb{F}^{1/4} S \mathbb{F}^{1/4}$ fulfill $\|U\| \leq \nu < 1$. Then

$$\begin{aligned} & \mathcal{K}(\mathbf{N}_{\bar{\mathbf{x}},\mathbb{F}} \parallel \mathbf{N}_{\bar{\mathbf{x}},\mathbb{Z}}) \\ &= -\log \det(\mathbb{I}_p + \mathbb{F}^{1/4} S \mathbb{F}^{1/4}) + \frac{1}{2} \text{tr} \{ \mathbb{F}(\mathbb{F}^{-1/2} + S)^2 - \mathbb{I}_p \} \\ &= -\log \det(\mathbb{I}_p + U) + \text{tr} U + \frac{1}{2} \text{tr}(\mathbb{F}S^2) \geq \frac{1}{2} \text{tr}(\mathbb{F}S^2). \end{aligned} \quad (1)$$

For two Gaussian distributions $\mathbf{N}_{\bar{x}, \mathbb{F}}, \mathbf{N}_{\bar{x}, \mathbb{Z}}$ with the same mean \bar{x}

$$\begin{aligned}\mathcal{K}(\mathbf{N}_{\bar{x}, \mathbb{F}} \parallel \mathbf{N}_{\bar{x}, \mathbb{Z}}) &= \frac{1}{2} \left\{ -\log \det(\mathbb{F} \mathbb{Z}^{-1}) + \text{tr}(\mathbb{F} \mathbb{Z}^{-1} - \mathbb{I}_p) \right\} \\ &= -\log \det \left\{ \mathbb{F}^{1/2} (\mathbb{F}^{-1/2} + S) \right\} + \frac{1}{2} \text{tr} \left\{ \mathbb{F} (\mathbb{F}^{-1/2} + S)^2 - \mathbb{I}_p \right\} \\ &= -\log \det(\mathbb{I}_p + U) + \frac{1}{2} \text{tr}(\mathbb{F} S^2 + 2\mathbb{F}^{1/2} S)\end{aligned}$$

and (1) follows by $x - \log(1 + x) \geq 0$ for any $x > -1$.

Consider symmetric matrices $S \in \mathfrak{M}_p$ such that for some $\nu < 1$

$$\|\mathbb{F}^{1/4} S \mathbb{F}^{1/4}\| \leq \nu. \quad (2)$$

Lemma

With $\gamma \sim \mathcal{N}(0, \mathbb{I}_p)$, $\mathbf{a} \in \mathbb{R}^p$, and $S \in \mathfrak{M}_p$ satisfying (2), define

$$H(\mathbf{a}, S) \stackrel{\text{def}}{=} -\log \det(\mathbb{F}^{-1/2} + S) - \mathbb{E}f(\bar{\mathbf{x}} + \mathbf{a} + (\mathbb{F}^{-1/2} + S)\gamma),$$

$$(\hat{\mathbf{a}}, \hat{S}) \stackrel{\text{def}}{=} \underset{(\mathbf{a}, S)}{\text{argmin}} H(\mathbf{a}, S).$$

Then the VI problem leads to minimization of the function $H(\mathbf{a}, S)$:

$$(\hat{\mathbf{x}}, \hat{\mathbb{Z}}) \stackrel{\text{def}}{=} \underset{(\mathbf{x}, \mathbb{Z})}{\text{argmin}} \mathcal{K}(\mathbb{P}_{\mathbf{x}, \mathbb{Z}} \parallel \mathbb{P}_f) = (\bar{\mathbf{x}} + \hat{\mathbf{a}}, (\mathbb{F}^{-1/2} + \hat{S})^{-2}).$$

1 Gaussian Variational Inference

2 **A basic lemma**

- Fourth order approximation
- Quadratic penalization

3 Solution to VI problem

4 Optimization vs sampling

Let $f(\mathbf{v})$ be a **smooth concave** function,

$$\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} f(\mathbf{v}), \quad \mathbb{F} = -\nabla^2 f(\mathbf{v}^*).$$

Let another function $g(\mathbf{v})$ satisfy for some vector \mathbf{A}

$$g(\mathbf{v}) - g(\mathbf{v}^*) = \langle \mathbf{v} - \mathbf{v}^*, \mathbf{A} \rangle + f(\mathbf{v}) - f(\mathbf{v}^*). \quad (3)$$

Define

$$\mathbf{v}^\circ \stackrel{\text{def}}{=} \operatorname{argmax}_{\mathbf{v}} g(\mathbf{v}), \quad g(\mathbf{v}^\circ) = \max_{\mathbf{v}} g(\mathbf{v}).$$

Aim: evaluate the quantities $\mathbf{v}^\circ - \mathbf{v}^*$ and $g(\mathbf{v}^\circ) - g(\mathbf{v}^*)$.

Let $L(\boldsymbol{v})$ be a **log-likelihood** function. Consider the MLE

$$\tilde{\boldsymbol{v}} = \operatorname{argmax}_{\boldsymbol{v}} L(\boldsymbol{v})$$

and the **background truth**

$$\boldsymbol{v}^* = \operatorname{argmax}_{\boldsymbol{v}} \mathbb{E}L(\boldsymbol{v})$$

Stochastically linear smooth (SLS) models: $\mathbb{E}L(\boldsymbol{v})$ is smooth in \boldsymbol{v}
and $\zeta(\boldsymbol{v}) = L(\boldsymbol{v}) - \mathbb{E}L(\boldsymbol{v})$ is linear in \boldsymbol{v} :

$$\boldsymbol{A} = \nabla\zeta(\boldsymbol{v}) = \nabla\zeta.$$

Let $h(\cdot)$ be concave and

$$\mathbf{v}^* = \operatorname{argmax} h(\mathbf{v}).$$

Consider

$$g(\mathbf{v}) = h(\mathbf{v}) - \|G\mathbf{v}\|^2/2,$$

$$f(\mathbf{v}) = h(\mathbf{v}) - \|G\mathbf{v}\|^2/2 + \langle G^2\mathbf{v}^*, \mathbf{v} \rangle, \quad .$$

Then $\nabla f(\mathbf{v}^*) = 0$ and $\mathbf{v}^* = \operatorname{argmax} f(\mathbf{v})$.

g is a linear perturbation of f with $\mathbf{A} = -G^2\mathbf{v}^*$.

Let f be a concave function and

$$\mathbf{v}^* = \operatorname{argmax} f(\mathbf{v}).$$

Let also \mathbf{v}° be a **current guess**. Define

$$g(\mathbf{v}) = f(\mathbf{v}) - \langle \nabla f(\mathbf{v}^\circ), \mathbf{v} - \mathbf{v}^\circ \rangle.$$

Then $\nabla g(\mathbf{v}^\circ) = 0$ and hence,

$$\mathbf{v}^\circ = \operatorname{argmax} g(\mathbf{v}).$$

g is a linear perturbation of f with $\mathbf{A} = \nabla f(\mathbf{v}^\circ)$.

Lemma

Let $f(\mathbf{v})$ be quadratic with $\nabla^2 f(\mathbf{v}) \equiv -\mathbb{F}$. If $g(\mathbf{v})$ satisfy (3), then

$$\mathbf{v}^\circ - \mathbf{v}^* = \mathbb{F}^{-1} \mathbf{A}, \quad g(\mathbf{v}^\circ) - g(\mathbf{v}^*) = \frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2.$$

Proof. Clearly $-\nabla^2 g(\mathbf{v}) \equiv -\mathbb{F}$ and

$$\nabla g(\mathbf{v}^*) - \nabla g(\mathbf{v}^\circ) = \mathbb{F}(\mathbf{v}^\circ - \mathbf{v}^*).$$

Further, (3) and $\nabla f(\mathbf{v}^*) = 0$ yield $\nabla g(\mathbf{v}^*) = \mathbf{A}$. Together with $\nabla g(\mathbf{v}^\circ) = 0$, this implies $\mathbf{v}^\circ - \mathbf{v}^* = \mathbb{F}^{-1} \mathbf{A}$.

Taylor expansion of g at \mathbf{v}° yields by $\nabla g(\mathbf{v}^\circ) = 0$

$$g(\mathbf{v}^*) - g(\mathbf{v}^\circ) = -\frac{1}{2} \|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{v}^*)\|^2 = -\frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2.$$

(\mathcal{T}_3^*) $f(\mathbf{v})$ is strongly concave, $\mathbb{D}^2(\mathbf{v}) \leq \nabla^2 f(\mathbf{v})$, and

$$\sup_{\mathbf{u}: \|\mathbb{D}(\mathbf{v})\mathbf{u}\| \leq r} \sup_{\mathbf{z} \in \mathbb{R}^p} \frac{|\langle \nabla^3 f(\mathbf{v} + \mathbf{u}), \mathbf{z}^{\otimes 3} \rangle|}{\|\mathbb{D}(\mathbf{v})\mathbf{z}\|^3} \leq \tau_3.$$

(\mathcal{T}_4^*) $f(\mathbf{v})$ is strongly concave, $\mathbb{D}^2(\mathbf{v}) \leq \nabla^2 f(\mathbf{v})$, and

$$\sup_{\mathbf{u}: \|\mathbb{D}(\mathbf{v})\mathbf{u}\| \leq r} \sup_{\mathbf{z} \in \mathbb{R}^p} \frac{|\langle \nabla^4 f(\mathbf{v} + \mathbf{u}), \mathbf{z}^{\otimes 4} \rangle|}{\|\mathbb{D}(\mathbf{v})\mathbf{z}\|^4} \leq \tau_4.$$

Banach's characterization [Banach, 1938] yields under (\mathcal{T}_3^*) (resp (\mathcal{T}_4^*))

$$|\langle \nabla^3 f(\mathbf{v} + \mathbf{u}), \mathbf{z}_1 \otimes \mathbf{z}_2 \otimes \mathbf{z}_3 \rangle| \leq \tau_3 \|\mathbb{D}(\mathbf{v})\mathbf{z}_1\| \|\mathbb{D}(\mathbf{v})\mathbf{z}_2\| \|\mathbb{D}(\mathbf{v})\mathbf{z}_3\|.$$

$$|\langle \nabla^4 f(\mathbf{v} + \mathbf{u}), \mathbf{z}_1 \otimes \mathbf{z}_2 \otimes \mathbf{z}_3 \otimes \mathbf{z}_4 \rangle| \leq \tau_4 \prod_{k=1}^4 \|\mathbb{D}(\mathbf{v})\mathbf{z}_k\|.$$

Proposition

Under (\mathcal{T}_3^*)

$$\begin{aligned} -\frac{2\tau_3}{3} \|\mathbb{F}^{-1/2} \mathbf{A}\|^3 &\leq 2g(\mathbf{v}^\circ) - 2g(\mathbf{v}^*) - \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 \\ &\leq \tau_3 \|\mathbb{F}^{-1/2} \mathbf{A}\|^3. \end{aligned} \quad (4)$$

and

$$\|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{v}^* - \mathbb{F}^{-1} \mathbf{A})\| \leq \frac{3\tau_3}{4} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2. \quad (5)$$

Implies Newton – Kantorovich – Nemirovskii-Nesterov results about quadratic convergence of second order methods.

By (\mathcal{T}_3^*) and $\nabla f(\mathbf{v}^*) = 0$, for any $\mathbf{v} \in \mathcal{A}(\mathbf{r})$

$$\begin{aligned} \left| f(\mathbf{v}^*) - f(\mathbf{v}) - \frac{1}{2} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\|^2 \right| &\leq \frac{\tau_3}{6} \|\mathbb{D}(\mathbf{v} - \mathbf{v}^*)\|^3 \\ &\leq \frac{\tau_3}{6} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\|^3. \end{aligned} \quad (6)$$

Further,

$$\begin{aligned} &g(\mathbf{v}) - g(\mathbf{v}^*) - \frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 \\ &= \langle \mathbf{v} - \mathbf{v}^*, \mathbf{A} \rangle + f(\mathbf{v}) - f(\mathbf{v}^*) - \frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 \\ &= -\frac{1}{2} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*) - \mathbb{F}^{-1/2} \mathbf{A}\|^2 + f(\mathbf{v}) - f(\mathbf{v}^*) + \frac{1}{2} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\|^2. \end{aligned}$$

As $\mathbf{v}^\circ \in \mathcal{A}(\mathbf{r})$ and it maximizes $g(\mathbf{v})$, we derive by (6) and Lemma 5

$$\begin{aligned}
 g(\mathbf{v}^\circ) - g(\mathbf{v}^*) - \frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 &= \max_{\mathbf{v} \in \mathcal{A}(\mathbf{r})} \left\{ g(\mathbf{v}) - g(\mathbf{v}^*) - \frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 \right\} \\
 &\leq \max_{\mathbf{v} \in \mathcal{A}(\mathbf{r})} \left\{ -\frac{1}{2} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*) - \mathbb{F}^{-1/2} \mathbf{A}\|^2 + \frac{\tau_3}{6} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\|^3 \right\} \\
 &\leq \frac{\tau_3}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^3.
 \end{aligned}$$

Now (4) follows from this and

$$\begin{aligned}
 g(\mathbf{v}^\circ) - g(\mathbf{v}^*) - \frac{1}{2} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 &\geq \max_{\mathbf{v} \in \mathcal{A}(\mathbf{r})} \left\{ -\frac{1}{2} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*) - \mathbb{F}^{-1/2} \mathbf{A}\|^2 - \frac{\tau_3}{6} \|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\|^3 \right\} \\
 &\geq -\frac{\tau_3}{3} \|\mathbb{F}^{-1/2} \mathbf{A}\|^3.
 \end{aligned}$$

For proving (5) use that $\nabla f(\mathbf{v}^*) = 0$, $\nabla g(\mathbf{v}^\circ) = 0$,
 $\nabla f(\mathbf{v}^\circ) = \nabla g(\mathbf{v}^\circ) - \mathbf{A} = -\mathbf{A}$, and $-\nabla^2 f(\mathbf{v}^*) = \mathbb{F}$. By Lemma ??
with $\mathbf{u} = \mathbb{F}^{-1}\mathbf{A}$

$$\|\mathbb{F}^{-1/2}\{\nabla f(\mathbf{v}^* + \mathbb{F}^{-1}\mathbf{A}) + \mathbf{A}\}\| \leq \frac{\tau_3}{2}\|\mathbb{F}^{-1/2}\mathbf{A}\|^2.$$

Further, by (3)

$$\begin{aligned}\|\mathbb{F}^{-1/2}\nabla g(\mathbf{v}^* + \mathbb{F}^{-1}\mathbf{A})\| &= \|\mathbb{F}^{-1/2}\{\nabla g(\mathbf{v}^* + \mathbb{F}^{-1}\mathbf{A}) - \mathbf{A} + \mathbf{A}\}\| \\ &\leq \|\mathbb{F}^{-1/2}\{\nabla f(\mathbf{v}^* + \mathbb{F}^{-1}\mathbf{A}) + \mathbf{A}\}\| \leq \frac{\tau_3}{2}\|\mathbb{F}^{-1/2}\mathbf{A}\|^2.\end{aligned}$$

By definition $\nabla g(\mathbf{v}^\circ) = 0$. This yields

$$\|\mathbb{F}^{-1/2}\{\nabla g(\mathbf{v}^* + \mathbb{F}^{-1}\mathbf{A}) - \nabla g(\mathbf{v}^\circ)\}\| \leq \frac{\tau_3}{2}\|\mathbb{F}^{-1/2}\mathbf{A}\|^2. \quad (7)$$

Now we can use with $\Delta = \mathbf{v}^* + \mathbb{F}^{-1} \mathbf{A} - \mathbf{v}^\circ$

$$\begin{aligned} & \mathbb{F}^{-1/2} \{ \nabla g(\mathbf{v}^* + \mathbb{F}^{-1} \mathbf{A}) - \nabla g(\mathbf{v}^\circ) \} \\ &= \left(\int_0^1 \mathbb{F}^{-1/2} \nabla^2 g(\mathbf{v}^\circ + t\Delta) \mathbb{F}^{-1/2} dt \right) \mathbb{F}^{1/2} \Delta. \end{aligned}$$

By (3) $\nabla^2 g(\mathbf{v}) = \nabla^2 f(\mathbf{v})$ for all \mathbf{v} . If $\|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}$, then (\mathcal{T}_3^*) implies $\|\mathbb{F}^{-1/2} \nabla^2 f(\mathbf{v}) \mathbb{F}^{-1/2} + \mathbb{I}_p\| \leq \omega^+ \leq \tau_3 \mathbf{r} \leq 1/3$. Hence,

$$\|\mathbb{F}^{-1/2} \{ \nabla g(\mathbf{v}^* + \mathbb{F}^{-1} \mathbf{A}) - \nabla g(\mathbf{v}^\circ) \}\| \geq (1 - \omega^+) \|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{v}^* - \mathbb{F}^{-1} \mathbf{A})\|.$$

This and (7) yield

$$\|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{v}^* - \mathbb{F}^{-1} \mathbf{A})\| \leq \frac{\tau_3}{2(1 - \omega^+)} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2 \leq \frac{3\tau_3}{4} \|\mathbb{F}^{-1/2} \mathbf{A}\|^2,$$

and (5) follows.

Lemma

For any $\xi \in \mathbb{R}^p$ with $\|\xi\| \leq 2r/3$ and τ with $\tau r \leq 1/2$, it holds

$$\max_{\|u\| \leq r} \left(\frac{\tau}{3} \|u\|^3 - \|u - \xi\|^2 \right) \leq \frac{\tau}{2} \|\xi\|^3, \quad (8)$$

$$\min_{\|u\| \leq r} \left(\frac{\tau}{3} \|u\|^3 + \|u - \xi\|^2 \right) \leq \frac{\tau}{3} \|\xi\|^3. \quad (9)$$

Any maximizer \mathbf{u} of the left hand-side of (8) satisfies

$$\tau \|\mathbf{u}\|^{1/2} \mathbf{u} - 2(\mathbf{u} - \boldsymbol{\xi}) = 0.$$

Therefore, $\mathbf{u} = \rho \boldsymbol{\xi}$ for some ρ , reducing the problem to the univariate case:

$$\max_{\|\mathbf{u}\| \leq r} \left(\frac{\tau}{3} \|\mathbf{u}\|^3 - \|\mathbf{u} - \boldsymbol{\xi}\|^2 \right) = \|\boldsymbol{\xi}\|^2 \max_{\rho: \|\rho \boldsymbol{\xi}\| \leq r} \left(\frac{\tau \|\boldsymbol{\xi}\|}{3} \rho^3 - (\rho - 1)^2 \right).$$

Define $a = \tau \|\boldsymbol{\xi}\|$. The conditions $\|\boldsymbol{\xi}\| \leq 2r/3$ and $\tau r \leq 1/2$ imply $a \leq 1/3$ and $\|\rho \boldsymbol{\xi}\| \leq r$ implies $|\rho| \leq 3/2$. The function $a\rho^3/3 - (\rho - 1)^2$ is concave on the interval $|\rho| \leq 3/2$ and hence, the maximizer ρ fulfills $a\rho^2 - 2\rho + 2 = 0$ yielding

$$\rho = \frac{1 \pm \sqrt{1 - 2a}}{a}, \quad |\rho| \leq 3/2.$$

As $a \in [0, 1/3]$, we can only use

$$\rho_a = \frac{1 - \sqrt{1 - 2a}}{a} = \frac{2}{1 + \sqrt{1 - 2a}}, \quad \rho_a - 1 = \frac{2a}{(1 + \sqrt{1 - 2a})^2}.$$

Therefore,

$$\begin{aligned} \max_{\|\mathbf{u}\| \leq r} \left(\frac{\tau}{3} \|\mathbf{u}\|^3 - \|\mathbf{u} - \boldsymbol{\xi}\|^2 \right) &= \frac{\tau \|\boldsymbol{\xi}\|^3 \rho_a^3}{3} - \|\boldsymbol{\xi}\|^2 (\rho_a - 1)^2 \\ &= \frac{\tau \|\boldsymbol{\xi}\|^3}{3} \frac{8(1 + \sqrt{1 - 2a}) - 12a}{(1 + \sqrt{1 - 2a})^4} \leq \frac{\tau \|\boldsymbol{\xi}\|^3}{3} \max_{a \in [0, 1/3]} \frac{8(1 + \sqrt{1 - 2a}) - 12a}{(1 + \sqrt{1 - 2a})^4} \end{aligned}$$

With $y = 1 + \sqrt{1 - 2a}$ or $-2a = (y - 1)^2 - 1 = y^2 - 2y$, represent

$$\phi(a) \stackrel{\text{def}}{=} \frac{8(1 + \sqrt{1 - 2a}) - 12a}{(1 + \sqrt{1 - 2a})^4} = \frac{8y + 6y^2 - 12y}{y^4} = \frac{6y - 4}{y^3},$$

and the latter decreases with $y \geq 1$. As $\phi(1/3) \leq 3/2$, (8) follows.

The proof of (9) is similar. The general case can be reduced to the univariate one by using $\mathbf{u} = \rho \boldsymbol{\xi}$. With $a = \tau \|\boldsymbol{\xi}\|$, the minimizer ρ_a reads as

$$\rho_a = \frac{2}{1 + \sqrt{1 + 2a}}, \quad 1 - \rho_a = \frac{\sqrt{1 + 2a} - 1}{\sqrt{1 + 2a} + 1} = \frac{2a}{(\sqrt{1 + 2a} + 1)^2},$$

yielding for $a \in [0, 1/3]$

$$\begin{aligned} \min_{\|\mathbf{u}\| \leq r} \left(\frac{\tau}{3} \|\mathbf{u}\|^3 + \|\mathbf{u} - \boldsymbol{\xi}\|^2 \right) &= \frac{\tau \|\boldsymbol{\xi}\|^3 \rho_a^3}{3} + \|\boldsymbol{\xi}\|^2 (\rho_a - 1)^2 \\ &\leq \frac{\tau \|\boldsymbol{\xi}\|^3}{3} \max_{a \in [0, 1/3]} \frac{8(1 + \sqrt{1 + 2a}) + 12a}{(1 + \sqrt{1 + 2a})^4}, \end{aligned}$$

and with $y = 1 + \sqrt{1 + 2a}$ or $2a = y^2 - 2y$,

$$\max_{a \in [0, 1/3]} \frac{8(1 + \sqrt{1 + 2a}) + 12a}{(1 + \sqrt{1 + 2a})^4} \leq \max_{y \geq 2} \frac{8y + 6y^2 - 12y}{y^4} = \max_{y \geq 2} \frac{6y - 4}{y^3} = 1$$

Proposition

Assume the conditions of Proposition 1 and (\mathcal{T}_4^*) . Then $\mathbf{v}^\circ = \operatorname{argmax}_{\mathbf{v}} g(\mathbf{v})$ satisfies $\|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{v}^*)\| \leq \mathbf{r}$. With

$$\mathbf{a} = \mathbb{F}^{-1}\{\mathbf{A} + \nabla\mathcal{T}(\mathbb{F}^{-1}\mathbf{A})\}$$

it holds

$$\|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{v}^* - \mathbf{a})\| \leq \frac{\tau_4 + 3\tau_3^2}{3} \|\mathbb{F}^{-1}\mathbf{A}\|^3.$$

Also with $\boldsymbol{\xi} = \mathbb{F}^{-1/2}\mathbf{A}$

$$\left| g(\mathbf{v}^\circ) - g(\mathbf{v}^*) - \frac{\|\boldsymbol{\xi}\|^2}{2} - \mathcal{T}(\mathbb{F}^{-1}\mathbf{A}) \right| \leq \frac{\tau_4 + 7\tau_3^2}{16} \|\boldsymbol{\xi}\|^4 + \frac{(\tau_4 + 3\tau_3^2)^2}{5} \|\boldsymbol{\xi}\|^6.$$

W.l.o.g. assume $\mathbf{v}^* = 0$. It holds by (\mathcal{T}_3^*)

$$\begin{aligned} \|\mathbb{F}^{1/2} \mathbf{a} - \boldsymbol{\xi}\| &= \|\mathbb{F}^{-1/2} \nabla \mathcal{T}(\mathbb{F}^{-1} \mathbf{A})\| \\ &= \sup_{\|\mathbf{u}\|=1} 3 \left| \langle \mathcal{T}, \mathbb{F}^{-1} \mathbf{A} \otimes \mathbb{F}^{-1} \mathbf{A} \otimes \mathbb{F}^{-1/2} \mathbf{u} \rangle \right| \leq \frac{\tau_3}{2} \|\boldsymbol{\xi}\|^2 \end{aligned} \quad (11)$$

yielding by $\|\boldsymbol{\xi}\| \leq \nu \mathbf{r}$

$$\|\mathbb{F}^{1/2} \mathbf{a}\| \leq \left(1 + \frac{\tau_3 \nu \mathbf{r}}{2} \right) \|\boldsymbol{\xi}\|. \quad (12)$$

Similarly for any \mathbf{v}

$$\|\mathbb{F}^{-1/2} \nabla^2 \mathcal{T}(\mathbb{F}^{-1/2} \mathbf{v}) \mathbb{F}^{-1/2}\| \leq \tau_3 \|\mathbf{v}\|.$$

Furthermore, the tensor $\nabla^2 \mathcal{T}(\mathbf{u})$ is linear in \mathbf{u} and hence,

$$\begin{aligned} & \sup_{t \in [0,1]} \|\mathbb{F}^{-1/2} \nabla^2 \mathcal{T}(t\mathbf{a} + (1-t)\mathbb{F}^{-1}\mathbf{A}) \mathbb{F}^{-1/2}\| \\ &= \max\{\|\mathbb{F}^{-1/2} \nabla^2 \mathcal{T}(\mathbb{F}^{-1}\mathbf{A}) \mathbb{F}^{-1/2}\|, \|\mathbb{F}^{-1} \nabla^2 \mathcal{T}(\mathbf{a})\|\} \\ &\leq \tau_3 \max\{\|\boldsymbol{\xi}\|, \|\mathbb{F}^{1/2}\mathbf{a}\|\}. \end{aligned}$$

Later we assume $\|\mathbb{F}^{1/2}\mathbf{a}\| \geq \|\mathbb{F}^{-1}\mathbf{A}\|$ in view of (12). This and (11) yield

$$\begin{aligned} & \|\mathbb{F}^{-1/2} \nabla \mathcal{T}(\mathbf{a}) - \mathbb{F}^{-1/2} \nabla \mathcal{T}(\mathbb{F}^{-1}\mathbf{A})\| \\ &\leq \sup_{t \in [0,1]} \|\mathbb{F}^{-1/2} \nabla^2 \mathcal{T}(t\mathbf{a} + (1-t)\mathbf{A}) \mathbb{F}^{-1/2}\| \|\mathbb{F}^{1/2}\mathbf{a} - \boldsymbol{\xi}\| \leq \frac{\tau_3^2}{2} \|\mathbb{F}^{1/2}\mathbf{a}\|^3 \end{aligned}$$

Further, in view of $\nabla \mathcal{T}(\mathbf{a}) = \frac{1}{2} \langle \nabla^3 f(\mathbf{v}^*), \mathbf{a} \otimes \mathbf{a} \rangle$

$$\|\mathbb{F}^{-1/2} \{\nabla f(\mathbf{a}) + \mathbb{F}\mathbf{a} - \nabla \mathcal{T}(\mathbf{a})\}\| \leq \frac{\tau_4}{6} \|\mathbb{F}^{1/2}\mathbf{a}\|^3.$$

Now we can bound the norm of $\mathbb{F}^{-1/2}\nabla g(\mathbf{a})$. In view of (3), it holds

$$\begin{aligned} \|\mathbb{F}^{-1/2}\nabla g(\mathbf{a})\| &= \|\mathbb{F}^{-1/2}\{\nabla g(\mathbf{a}) + \mathbb{F}\mathbf{a} - \nabla\mathcal{T}(\mathbf{A}) - \mathbf{A}\}\| \\ &\leq \|\mathbb{F}^{-1/2}\{\nabla f(\mathbf{a}) + \mathbb{F}\mathbf{a} - \nabla\mathcal{T}(\mathbf{a})\}\| + \|\mathbb{F}^{-1/2}\{\nabla\mathcal{T}(\mathbf{a}) - \nabla\mathcal{T}(\mathbf{A})\}\| \\ &\leq \frac{\tau_4 + 3\tau_3^2}{6} \|\mathbb{F}^{1/2}\mathbf{a}\|^3. \end{aligned}$$

By definition $\nabla g(\mathbf{v}^\circ) = 0$. This yields

$$\|\mathbb{F}^{-1/2}\{\nabla g(\mathbf{a}) - \nabla g(\mathbf{v}^\circ)\}\| \leq \frac{\tau_4 + 3\tau_3^2}{6} \|\mathbb{F}^{1/2}\mathbf{a}\|^3. \quad (13)$$

Furthermore, with $\Delta = \mathbf{a} - \mathbf{v}^\circ$

$$\mathbb{F}^{-1/2}\{\nabla g(\mathbf{a}) - \nabla g(\mathbf{v}^\circ)\} = \left(\int_0^1 \mathbb{F}^{-1/2} \nabla^2 g(\mathbf{v}^\circ + t\Delta) \mathbb{F}^{-1/2} dt \right) \mathbb{F}^{1/2} \Delta.$$

By (3) $\nabla^2 g(\mathbf{v}) = \nabla^2 f(\mathbf{v})$ for all \mathbf{v} . If $\|\mathbb{F}^{1/2}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}$, then (\mathcal{T}_3^*) implies $\|\mathbb{F}^{-1/2} \nabla^2 f(\mathbf{v}) \mathbb{F}^{-1/2} + \mathbb{I}_p\| \leq \omega^+$ with $\omega^+ \leq \tau_3 \mathbf{r} \leq 1/3$ and

$$\|\mathbb{F}^{-1/2}\{\nabla g(\mathbf{a}) - \nabla g(\mathbf{v}^\circ)\}\| \geq (1 - \tau_3 \mathbf{r})\|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{a})\|.$$

This, (12), and (13) yield in view of $\tau_3 \mathbf{r} \leq 1/3$ and $\nu = 2/3$

$$\|\mathbb{F}^{1/2}(\mathbf{v}^\circ - \mathbf{a})\| \leq \frac{\tau_4 + 3\tau_3^2}{6(1 - \tau_3 \mathbf{r})} \|\mathbb{F}^{1/2} \mathbf{a}\|^3 \leq \frac{\tau_4 + 3\tau_3^2}{3} \|\boldsymbol{\xi}\|^3. \quad (14)$$

It remains to bound $g(\mathbf{v}^\circ) - g(0)$. By (11)

$$\frac{1}{2}\|\boldsymbol{\xi}\|^2 - \langle \mathbf{A}, \mathbf{a} \rangle + \frac{1}{2}\|\mathbb{F}^{1/2}\mathbf{a}\|^2 = \frac{1}{2}\|\mathbb{F}^{1/2}\mathbf{a} - \boldsymbol{\xi}\|^2 \leq \frac{\tau_3^2}{8}\|\boldsymbol{\xi}\|^4.$$

First consider $g(\mathbf{a}) - g(0)$. One more use of (\mathcal{T}_4^*) yields with $\mathbf{v}^* = 0$ and $-\nabla^2 f(0) = \mathbb{F}$

$$\begin{aligned} & \left| g(\mathbf{a}) - g(0) - \frac{1}{2}\|\boldsymbol{\xi}\|^2 - \mathcal{T}(\mathbf{a}) \right| \\ &= \left| f(\mathbf{a}) - f(0) + \langle \mathbf{A}, \mathbf{a} \rangle - \frac{1}{2}\|\boldsymbol{\xi}\|^2 - \mathcal{T}(\mathbf{a}) \right| \\ &\leq \left| f(\mathbf{a}) - f(0) + \frac{1}{2}\|\mathbb{F}^{1/2}\mathbf{a}\|^2 - \mathcal{T}(\mathbf{a}) \right| + \frac{\tau_3^2}{8}\|\boldsymbol{\xi}\|^4 \\ &\leq \frac{\tau_4}{24}\|\mathbb{F}^{1/2}\mathbf{a}\|^4 + \frac{\tau_3^2}{8}\|\boldsymbol{\xi}\|^4 \leq \frac{\tau_4 + 2\tau_3^2}{16}\|\boldsymbol{\xi}\|^4. \end{aligned}$$

Also by $\nabla g(\mathbf{v}^\circ) = 0$ and (14), it holds for some $\mathbf{v} \in [\mathbf{a}, \mathbf{v}^\circ]$ as in (14)

$$\begin{aligned} |g(\mathbf{a}) - g(\mathbf{v}^\circ)| &\leq \frac{1}{2} \|\mathbb{F}^{-1/2} \nabla^2 g(\mathbf{v}) \mathbb{F}^{-1/2}\| \|\mathbb{F}^{1/2}(\mathbf{a} - \mathbf{v}^\circ)\|^2 \\ &\leq \frac{(\tau_4 + 3\tau_3^2)^2}{72(1 - \tau_3 \mathbf{r})^3} \|\mathbb{F}^{1/2} \mathbf{a}\|^6 < \frac{(\tau_4 + 3\tau_3^2)^2}{5} \|\boldsymbol{\xi}\|^6, \end{aligned}$$

Moreover, similarly to (11)

$$\begin{aligned} |\mathcal{T}(\mathbf{a}) - \mathcal{T}(\mathbb{F}^{-1} \mathbf{A})| &\leq \sup_{t \in [0,1]} \|\mathbb{F}^{-1/2} \nabla \mathcal{T}(t\mathbb{F}^{-1} \mathbf{A} + (1-t)\mathbf{a})\| \|\mathbb{F}^{1/2} \mathbf{a} - \boldsymbol{\xi}\| \\ &\leq \frac{\tau_3^2}{4} \|\mathbb{F}^{1/2} \mathbf{a}\|^2 \|\boldsymbol{\xi}\|^2 \leq \frac{5\tau_3^2}{16} \|\boldsymbol{\xi}\|^4. \end{aligned}$$

Summing up the obtained bounds yields (10).

Here we discuss the case when $g(\mathbf{v}) - f(\mathbf{v})$ is quadratic. The general case can be reduced to the situation with $g(\mathbf{v}) = f(\mathbf{v}) - \|G\mathbf{v}\|^2/2$. To make the dependence of G more explicit, denote $f_G(\mathbf{v}) \stackrel{\text{def}}{=} f(\mathbf{v}) - \|G\mathbf{v}\|^2/2$,

$$\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} f(\mathbf{v}),$$

$$\mathbf{v}_G^* = \operatorname{argmax}_{\mathbf{v}} f_G(\mathbf{v}) = \operatorname{argmax}_{\mathbf{v}} \{f(\mathbf{v}) - \|G\mathbf{v}\|^2/2\}.$$

We study the bias $\mathbf{v}_G^* - \mathbf{v}^*$ induced by this penalization.

Lemma

Let $f(\mathbf{v})$ be quadratic with $\mathbb{F} \equiv -\nabla^2 f(\mathbf{v})$ and $\mathbf{S}_G \equiv G^2 \mathbf{v}^*$. Then it holds with $\mathbb{F}_G = \mathbb{F} + G^2$

$$\begin{aligned}\mathbf{v}^* - \mathbf{v}_G^* &= \mathbb{F}_G^{-1} \mathbf{S}_G = -\mathbb{F}_G^{-1} G^2 \mathbf{v}^*, \\ f_G(\mathbf{v}_G^*) - f_G(\mathbf{v}^*) &= \frac{1}{2} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^2 = \frac{1}{2} \|\mathbb{F}_G^{-1/2} G^2 \mathbf{v}^*\|^2.\end{aligned}$$

Quadraticity of $f(\mathbf{v})$ implies quadraticity of $f_G(\mathbf{v})$ with $\nabla^2 f_G(\mathbf{v}) \equiv -\mathbb{F}_G$ and

$$\nabla f_G(\mathbf{v}^*) - \nabla f_G(\mathbf{v}_G^*) = \mathbb{F}_G (\mathbf{v}_G^* - \mathbf{v}^*).$$

Further, $\nabla f(\mathbf{v}^*) = 0$ yielding $\nabla f_G(\mathbf{v}^*) = \mathbf{S}_G = G^2 \mathbf{v}^*$. Together with $\nabla f_G(\mathbf{v}_G^*) = 0$, this implies $\mathbf{v}^* - \mathbf{v}_G^* = \mathbb{F}_G^{-1} \mathbf{S}_G$. The Taylor expansion of f_G at \mathbf{v}_G^* yields

$$f_G(\mathbf{v}^*) - f_G(\mathbf{v}_G^*) = -\frac{1}{2} \|\mathbb{F}_G^{1/2} (\mathbf{v}^* - \mathbf{v}_G^*)\|^2 = -\frac{1}{2} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^2$$

and the assertion follows.

Proposition

Let f be concave and $\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} f(\mathbf{v})$. Define $\mathbf{S}_G = G^2 \mathbf{v}^*$,

$$\mathbb{F}_G = -\nabla^2 f(\mathbf{v}^*) + G^2, \quad \mathbf{b}_G = \|\mathbb{F}_G^{-1/2} G^2 \mathbf{v}^*\| = \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|.$$

With $\nu = 2/3$, assume (\mathcal{T}_3^*) for $\mathbf{r} = \nu^{-1} \mathbf{b}_G$ and $\mathbb{D}^2 \leq \mathbb{F}_G$. Then $\|\mathbb{F}_G^{1/2}(\mathbf{v}_G^* - \mathbf{v}^*)\| \leq \nu^{-1} \mathbf{b}_G$ or, equivalently,

$$\mathbf{v}_G^* \in \mathcal{A}_G \stackrel{\text{def}}{=} \{\mathbf{v}: \|\mathbb{F}_G^{1/2}(\mathbf{v} - \mathbf{v}^*)\| \leq \nu^{-1} \mathbf{b}_G\}. \quad (15)$$

Moreover,

$$\begin{aligned} \|\mathbb{F}_G^{1/2}(\mathbf{v}^* - \mathbf{v}_G^*) - \mathbb{F}_G^{-1/2} \mathbf{S}_G\|^2 &\leq \tau_3 \mathbf{b}_G^3, \\ |2f_G(\mathbf{v}_G^*) - 2f_G(\mathbf{v}^*) - \mathbf{b}_G^2| &\leq \tau_3 \mathbf{b}_G^3. \end{aligned}$$

With $\mathbf{S}_G = G^2 \mathbf{v}^*$, define $g_G(\mathbf{v})$ by

$$g_G(\mathbf{v}) - g_G(\mathbf{v}_G^*) = f_G(\mathbf{v}) - f_G(\mathbf{v}_G^*) + \langle \mathbf{S}_G, \mathbf{v} - \mathbf{v}_G^* \rangle. \quad (16)$$

The function f_G is concave, the same holds for g_G from (16). Now we show that $\mathbf{v}^* = \operatorname{argmax} g_G(\mathbf{v})$. It suffices to check that $\nabla g_G(\mathbf{v}^*) = 0$. Indeed, by definition, $\nabla f(\mathbf{v}^*) = 0$, and hence, $\nabla f_G(\mathbf{v}^*) = -G^2 \mathbf{v}^* + \mathbf{S}_G = 0$. Now the results follow from Proposition 1 applied with $f(\mathbf{v}) = g_G(\mathbf{v}) = f_G(\mathbf{v}) - \langle \mathbf{A}, \mathbf{v} \rangle$, $g(\mathbf{v}) = f_G(\mathbf{v})$, and $\mathbf{A} = \mathbf{S}_G$.

Define $\mathbb{F}_G = -\nabla^2 f(\mathbf{v}^*) + G^2$, $\mathbf{S}_G = G^2 \mathbf{v}^*$, and

$$\mathbf{m}_G = \mathbb{F}_G^{-1} \{ \mathbf{S}_G + \nabla \mathcal{T}(\mathbb{F}_G^{-1} \mathbf{S}_G) \}$$

with $\mathcal{T}(\mathbf{u}) = \frac{1}{6} \langle \nabla^3 f(\mathbf{v}^*), \mathbf{u}^{\otimes 3} \rangle$.

(\mathcal{T}_4^*) $f(\mathbf{v})$ is strongly concave, $\mathbb{D}^2(\mathbf{v}) \leq -\nabla^2 f(\mathbf{v})$, and

$$\sup_{\mathbf{u}: \|\mathbb{D}(\mathbf{v})\mathbf{u}\| \leq \mathbf{r}} \sup_{\mathbf{z} \in \mathbb{R}^p} \frac{|\langle \nabla^4 f(\mathbf{v} + \mathbf{u}), \mathbf{z}^{\otimes 4} \rangle|}{\|\mathbb{D}(\mathbf{v})\mathbf{z}\|^4} \leq \tau_4.$$

Typically $\tau_3 \asymp n^{-1/2}$ and $\tau_4 \asymp n^{-1}$.

Proposition

Let f be concave and $\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} f(\mathbf{v})$. With $\nu = 2/3$, assume (\mathcal{T}_3^*) and (\mathcal{T}_4^*) for $\mathbf{r} = \mathbf{r}_G \stackrel{\text{def}}{=} \nu^{-1} \mathbf{b}_G$ and $\mathbb{D}^2 \leq \mathbb{F}_G$. Then (15) holds. Furthermore, $\|\mathbb{F}_G^{1/2} \mathbf{m}_G\| \leq \mathbf{r}_G$ and

$$\|\mathbb{F}_G^{1/2} \mathbf{m}_G - \mathbb{F}_G^{-1/2} \mathbf{S}_G\| \leq \frac{\tau_3}{2} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^2 \leq \frac{\tau_3 \nu \mathbf{r}_G}{2} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|,$$

$$\|\mathbb{F}_G^{1/2} \mathbf{m}_G\| \leq \left(1 + \frac{\tau_3 \nu \mathbf{r}_G}{2}\right) \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|,$$

$$\|\mathbb{F}_G^{1/2} (\mathbf{v}^* - \mathbf{v}_G^* - \mathbf{m}_G)\| \leq \frac{\tau_4 + 3\tau_3^2}{3} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^3.$$

Also

$$\begin{aligned} & \left| f_G(\mathbf{v}_G^*) - f_G(\mathbf{v}^*) - \frac{1}{2} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^2 - \mathcal{T}(\mathbf{m}_G) \right| \\ & \leq \frac{\tau_4 + 2\tau_3^2}{16} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^4 + \frac{(\tau_4 + 3\tau_3^2)^2}{5} \|\mathbb{F}_G^{-1/2} \mathbf{S}_G\|^6. \end{aligned}$$

1 Gaussian Variational Inference

2 A basic lemma

- Fourth order approximation
- Quadratic penalization

3 Solution to VI problem

4 Optimization vs sampling

For $X \sim \mathbb{P}_f \propto e^{f(\mathbf{x})}$, consider

$$\bar{\mathbf{x}} = \mathbb{E}_f X, \quad \Sigma = \text{Var}(X), \quad \mathbb{F} = -\nabla^2 f(\bar{\mathbf{x}}).$$

Consider

$$H(\mathbf{a}, S) \stackrel{\text{def}}{=} -\log \det(\mathbb{F}^{-1/2} + S) - \mathbb{E} f(\bar{\mathbf{x}} + \mathbf{a} + (\mathbb{F}^{-1/2} + S)\boldsymbol{\gamma}),$$
$$(\hat{\mathbf{a}}, \hat{S}) \stackrel{\text{def}}{=} \underset{(\mathbf{a}, S)}{\text{argmin}} H(\mathbf{a}, S).$$

A guess $(\mathbf{a}, S) = (0, 0)$. How far from the solution $(\hat{\mathbf{a}}, \hat{S})$?

Technical issue: **anisotropic** smoothness in \mathbf{a} and S directions.

Fix $\mathbb{Z}^{-1/2} = \mathbb{F}^{-1/2} + S$ and optimize w.r.t. \mathbf{a} .

For $\mathbb{Z}^{-1/2} = \mathbb{F}^{-1/2} + S$ fixed, consider $H(\mathbf{a}) = H(\mathbf{a}, S)$

$$\hat{\mathbf{a}} \stackrel{\text{def}}{=} \underset{\mathbf{a}}{\operatorname{argmin}} H(\mathbf{a}) = \underset{\mathbf{a}}{\operatorname{argmax}} \mathbb{E} f(\bar{\mathbf{x}} + \mathbf{a} + \mathbb{Z}^{-1/2} \boldsymbol{\gamma}).$$

Main step: compute $\mathbf{A} = \nabla H(0)$ and $\mathcal{F} = -\nabla^2 H(0)$.

A guess: $\mathcal{F} \approx \mathbb{F} = -\nabla^2 f(\bar{\mathbf{x}})$, $\mathbf{A} \approx 0$ up to fourth order.

Fix \mathbf{a} and consider

$$h(t) = -\mathbb{E}f(\bar{\mathbf{x}} + t\mathbf{a} + \mathbb{Z}^{-1/2}\boldsymbol{\gamma}).$$

Lemma

The function $h(t) = H(t\mathbf{a})$ is strongly convex and satisfies

$$h''(t) = -\langle \mathbb{E}\nabla^2 f(\bar{\mathbf{x}} + t\mathbf{a} + \mathbb{Z}^{-1/2}\boldsymbol{\gamma}), \mathbf{a}^{\otimes 2} \rangle.$$

Concavity of $f(\cdot)$ implies convexity of h .

Lemma

It holds with $\mathbb{F} = -\nabla^2 f(\bar{\mathbf{x}})$

$$h''(0) = -\mathbb{E} \langle \nabla^2 f(\bar{\mathbf{x}} + \gamma_{\mathbb{Z}}), \mathbf{a}^{\otimes 2} \rangle,$$

and with $p = \text{tr}(\mathbb{D} \mathbb{F}^{-1} \mathbb{D})$ and $\alpha = \|\mathbb{D} \mathbb{F}^{-1} \mathbb{D}\|$

$$|h''(0) - \mathbf{a}^\top \mathbb{F} \mathbf{a}| \leq \frac{\tau_4(p + 2\alpha)}{2} \|\mathbb{D} \mathbf{a}\|^2. \quad (17)$$

It holds

$$-\langle \nabla^2 f(\bar{\mathbf{x}}), \mathbf{a}^{\otimes 2} \rangle = \mathbf{a}^\top \mathbb{F} \mathbf{a}.$$

For any $\mathbf{u} \in \mathbb{R}^p$,

$$\begin{aligned} & |-\langle \nabla^2 f(\bar{\mathbf{x}} + \gamma_{\mathbf{z}}), \mathbf{u}^{\otimes 2} \rangle + \langle \nabla^2 f(\bar{\mathbf{x}}), \mathbf{u}^{\otimes 2} \rangle + \langle \nabla^3 f(\bar{\mathbf{x}}), \gamma_{\mathbf{z}} \otimes \mathbf{u}^{\otimes 2} \rangle| \\ & \leq \frac{1}{2} \tau_4 \|\mathbb{D} \gamma_{\mathbf{z}}\|^2 \|\mathbb{D} \mathbf{u}\|^2. \end{aligned}$$

With $p = \text{tr}(\mathbb{D}^2 \mathbb{F}^{-1})$

$$\mathbb{E} \|\mathbb{D} \gamma_{\mathbf{z}}\|^2 = p.$$

Further, $\mathbb{E} \langle \nabla^3 f(\bar{\mathbf{x}}), \gamma_{\mathbf{z}} \otimes \mathbf{a}^{\otimes 2} \rangle = 0$ and (17) follows.

Define for any direction \mathbf{a}

$$h(t) = -\mathbb{E}f(\bar{\mathbf{x}} + t\mathbf{a} + \mathbb{Z}^{-1/2}\boldsymbol{\gamma}).$$

Lemma

It holds with $p = \text{tr}(\mathbb{D}\mathbb{F}^{-1}\mathbb{D})$, $\alpha = \|\mathbb{D}\mathbb{F}^{-1}\mathbb{D}\|$,

$$|h'(0)| \leq \frac{\tau_4 (p + \alpha)^{3/2}}{6} \|\mathbb{D}\mathbf{a}\| + \frac{\diamond_{4,1}}{1 - \diamond} \|\mathbb{D}\mathbf{a}\|.$$

With $\gamma_Z = Z^{-1/2}\gamma$, Taylor expansion of $\nabla f(\bar{\mathbf{x}} + \gamma_Z)$ yields for any $\mathbf{u} \in \mathbb{R}^p$

$$\begin{aligned}
 & |\langle \nabla f(\bar{\mathbf{x}} + \gamma_Z), \mathbf{u} \rangle - \langle \nabla f(\bar{\mathbf{x}}), \mathbf{u} \rangle - \langle \nabla^2 f(\bar{\mathbf{x}}), \gamma_Z \otimes \mathbf{u} \rangle \\
 & \quad - \frac{1}{2} \langle \nabla^3 f(\bar{\mathbf{x}}), \gamma_F \otimes \gamma_Z \otimes \mathbf{u} \rangle| \leq \frac{1}{6} \tau_4 \|\mathbb{D}\gamma_Z\|^3 \|\mathbb{D}\mathbf{u}\|. \quad (18)
 \end{aligned}$$

Also by Laplace approximation

$$|\langle \nabla f(\bar{\mathbf{x}}), \mathbf{a} \rangle - \frac{1}{2} \mathbb{E} \langle \nabla^3 f(\bar{\mathbf{x}}), \gamma_F \otimes \gamma_F \otimes \mathbf{a} \rangle| \leq \frac{\diamond_{4,1}}{1 - \diamond} \|\mathbb{D}\mathbf{a}\|.$$

Now we apply (18) with $\mathbf{u} = \mathbf{a}$ and $\mathbb{E} \|\mathbb{D}\gamma_F\|^3 \leq (p + \alpha)^{3/2}$. The use of $\mathbb{E} \langle \nabla^2 f(\bar{\mathbf{x}}), \gamma_F \otimes \mathbf{a} \rangle = 0$ yields

$$|\mathbb{E} \langle \nabla f(\bar{\mathbf{x}} + Z^{-1/2}\gamma), \mathbf{a} \rangle| \leq \frac{\tau_4 (p + \alpha)^{3/2}}{6} \|\mathbb{D}\mathbf{a}\| + \frac{\diamond_{4,1}}{1 - \diamond} \|\mathbb{D}\mathbf{a}\|.$$

Theorem (3-bound)

$$\|\mathbb{F}^{1/2}\hat{\mathbf{a}} - \mathbb{F}^{-1/2}\mathbf{A}\| \leq \tau_3 \|\mathbb{F}^{-1/2}\mathbf{A}\|^3$$

Theorem (4-bound)

$$\|\mathbb{F}^{1/2}\hat{\mathbf{a}} - \mathbb{F}^{-1/2}\mathbf{A} - \mathbb{F}^{-1/2}\nabla\mathcal{T}(\mathbb{F}^{-1}\mathbf{A})\| \leq \mathbf{C}(\tau_3^2 + \tau_4) \|\mathbb{F}^{-1/2}\mathbf{A}\|^3.$$

1 Gaussian Variational Inference

2 A basic lemma

- Fourth order approximation
- Quadratic penalization

3 Solution to VI problem

4 Optimization vs sampling

Optimization: $f(\mathbf{x}) \rightarrow \min$.

Sampling: $\mathbf{X} \propto \exp\{-f(\mathbf{x}) + \log \pi(\mathbf{x})\}$ for a sampling density π .

Sampling gradient free procedure (1 step):

- draw a mini-batch $\mathbf{X}_1, \dots, \mathbf{X}_n$ from π ;
- compute $w_i = e^{-f(\mathbf{X}_i)}$;
- mini-batch averaging

$$\bar{\mathbf{x}}_\pi = \frac{\sum_{i=1}^n \mathbf{X}_i w_i}{\sum_{i=1}^n w_i}.$$

Update of π (EnKF, diffusion models, VAE, etc.), loop.

Issues: mini-batch size n , averaging over steps, rate $\|\bar{\mathbf{x}}_\pi - \bar{\mathbf{x}}\|$.



Alquier, P. and Ridgway, J. (2020).
Concentration of tempered posteriors and of their variational approximations.
The Annals of Statistics, 48(3):1475 – 1497.



Banach, S. (1938).
Über homogene Polynome in (L^2) .
Studia Mathematica, 7(1):36–44.



Challis, E. and Barber, D. (2013).
Gaussian kullback-leibler approximate inference.
Journal of Machine Learning Research, 14(68):2239–2286.



David M. Blei, A. K. and McAuliffe, J. D. (2017).
Variational inference: A review for statisticians.
Journal of the American Statistical Association, 112(518):859–877.



Han, W. and Yang, Y. (2019).
Statistical inference in mean-field variational bayes.
<https://arxiv.org/abs/1911.01525>.



Katsevich, A. and Rigollet, P. (2023).
On the approximation accuracy of gaussian variational inference.



Lambert, M., Chewi, S., Bach, F., Bonnabel, S., and Rigollet, P. (2023).
Variational inference via wasserstein gradient flows.
<https://arxiv.org/abs/2205.15902>.



Wang, Y. and Blei, D. M. (2019).
Frequentist consistency of variational bayes.
Journal of the American Statistical Association, 114(527):1147–1161.



Zhang, F. and Gao, C. (2020).
Convergence rates of variational posterior distributions.
The Annals of Statistics, 48(4):2180 – 2207.