# Calculating Optimistic Likelihoods Using (Geodesically) Convex Optimization

Viet Anh Nguyen[†], Soroosh Shafieezadeh-Abadeh[†], Man-Chung Yue[‡],
Daniel Kuhn[†] and Wolfram Wiesemann[§]

[†]EPFL, [‡]The Hong Kong Polytechnic University, [§]Imperial College London

## Contributions

MLE problem for i.i.d. data points
$$x_1^M \triangleq x_1, \ldots, x_M \in \mathbb{R}^n$$
and candidate distributions $P_c = \mathcal{N}(\mu_c, \Sigma_c)$:
$$c^\star \in \arg\max_{c \in \mathcal{C}} \left\{ \ell(x_1^M, P_c) \triangleq -\frac{1}{M}\sum_{m=1}^M (x_m - \mu_c)^\top \Sigma_c^{-1}(x_m - \mu_c) - \log\det \Sigma_c \right\}$$

**Motivation:**
- MLE problem is fundamental in hypothesis testing and discriminant analysis
- The parameters $(\mu_c, \Sigma_c)$ of $P_c$ are uncertain
- Ignoring this uncertainty leads to poor out-of-sample performance

**Contributions:**
- We propose an **optimistic likelihood** (OL) problem over **Fisher-Rao** (FR) and **Kullback-Leibler** (KL) ambiguity sets containing normal distributions
- For FR ambiguity sets, the OL problem reduces to a **geodesically convex** problem
- We devise a **Riemannian gradient descent** algorithm
- For KL ambiguity sets, the OL problem reduces to a one dimensional convex problem

## Fisher-Rao Distance

Parametric distributions with density function $p_\theta(x)$
- For any $\theta \in \Theta$, the Fisher information matrix
$$I_\theta \triangleq \mathbb{E}_x[\nabla_\theta \log(p_\theta(x)) \nabla_\theta \log(p_\theta(x))^\top]$$
defines an inner product $\langle \cdot, \cdot \rangle_\theta$ on the tangent space $T_\theta \Theta$ as
$$\langle \zeta_1, \zeta_2 \rangle_\theta = \zeta_1^T I_\theta \zeta_2, \quad \forall \zeta_1, \zeta_2 \in T_\theta\Theta$$
- The set of $\{\langle \cdot, \cdot \rangle_\theta\}_{\theta \in \Theta}$ defines a Riemannian metric called the FR metric
- The FR metric is invariant under transformations on the data space
- The FR distance on $\Theta$ is a geodesic distance defined as
$$d(\theta_0, \theta_1) = \inf_\gamma \int_0^1 \sqrt{\langle \gamma'(t), \gamma'(t) \rangle_{\gamma(t)}} dt$$
- The infimum is over smooth curves $\gamma : [0,1] \to \Theta$ with $\gamma(0) = \theta_0$ and $\gamma(1) = \theta_1$

**Proposition 1.** [Atkinson and Mitchell (1981)] For the family of Gaussian distributions with identical mean and and covariance matrices $\Sigma_0, \Sigma_1$, we have
$$d(\Sigma_0, \Sigma_1) = \frac{1}{\sqrt{2}} \left\| \log\left(\Sigma_1^{-\frac{1}{2}} \Sigma_0 \Sigma_1^{-\frac{1}{2}}\right) \right\|_F$$

## Kullback-Leibler Divergence

For distributions $P_0$ and $P_1$ with density functions $p_0(x)$ and $p_1(x)$, we have
$$\text{KL}(P_0 \| P_1) = \int_{-\infty}^{\infty} p_0(x) \log\left(\frac{p_0(x)}{q_1(x)}\right) dx$$
When $P_0 = \mathcal{N}(\mu_0, \Sigma_0)$ and $P_1 = \mathcal{N}(\mu_1, \Sigma_1)$, the KL divergence coincides with
$$\text{KL}(P_0 \| P_1) = \frac{1}{2} \left( \text{Tr}\left[\Sigma_1^{-1}\Sigma_0\right] + \log\det(\Sigma_1\Sigma_0^{-1}) - n + (\mu_0 - \mu_1)^\top \Sigma_1^{-1}(\mu_0 - \mu_1) \right)$$

## Optimistic Likelihood Problems

We consider the optimistic likelihood problem
$$\text{OL:} \quad \max_{P \in \mathcal{P}} \ell(x_1^M, P) \quad \text{with} \quad \mathcal{P} = \left\{ P \in \mathcal{M} : \varphi(\hat{P}, P) \leq \rho \right\}$$

- Candidate Gaussian distribution: $\hat{P} = \mathcal{N}(\hat\mu, \hat\Sigma)$
- $\mathcal{M}$ is the family of Gaussian distributions with fixed mean $\hat\mu$
- $\varphi(\cdot, \cdot)$ is the dissimilarity measure $\Rightarrow$ FR distance or KL divergence
- $\rho$ is the size of the ambiguity set

## OL Problem under the FR Distance

The **OL** problem reduces to $\min_{\Sigma \in \mathcal{B}^{\text{FR}}} L(\Sigma)$, where
$$L(\Sigma) \triangleq \langle S, \Sigma^{-1}\rangle + \log\det \Sigma$$
$$\mathcal{B}^{\text{FR}} \triangleq \{\Sigma \in \mathbb{S}_{++}^n : d(\Sigma, \hat\Sigma) \leq \rho\}$$
$$S = M^{-1}\sum_{m=1}^M (x_m - \hat\mu)(x_m - \hat\mu)^\top$$

**Theorem 1.** ▶ $\mathcal{B}^{\text{FR}}$ is a geodesically convex set
- $L(\cdot)$ is a geodesically convex function over $\mathbb{S}_{++}^n$
- $L(\cdot)$ is geodesically $\beta$-smooth and $\sigma$-strongly on $\mathcal{B}^{\text{FR}}$ with
$$\beta = \frac{2\lambda_{\max}(S)}{\lambda_{\min}(\hat\Sigma)\exp(-\sqrt{2}\rho)}, \quad \text{and} \quad \sigma = \frac{2\lambda_{\min}(S)}{\lambda_{\max}(\hat\Sigma)\exp(\sqrt{2}\rho)}.$$
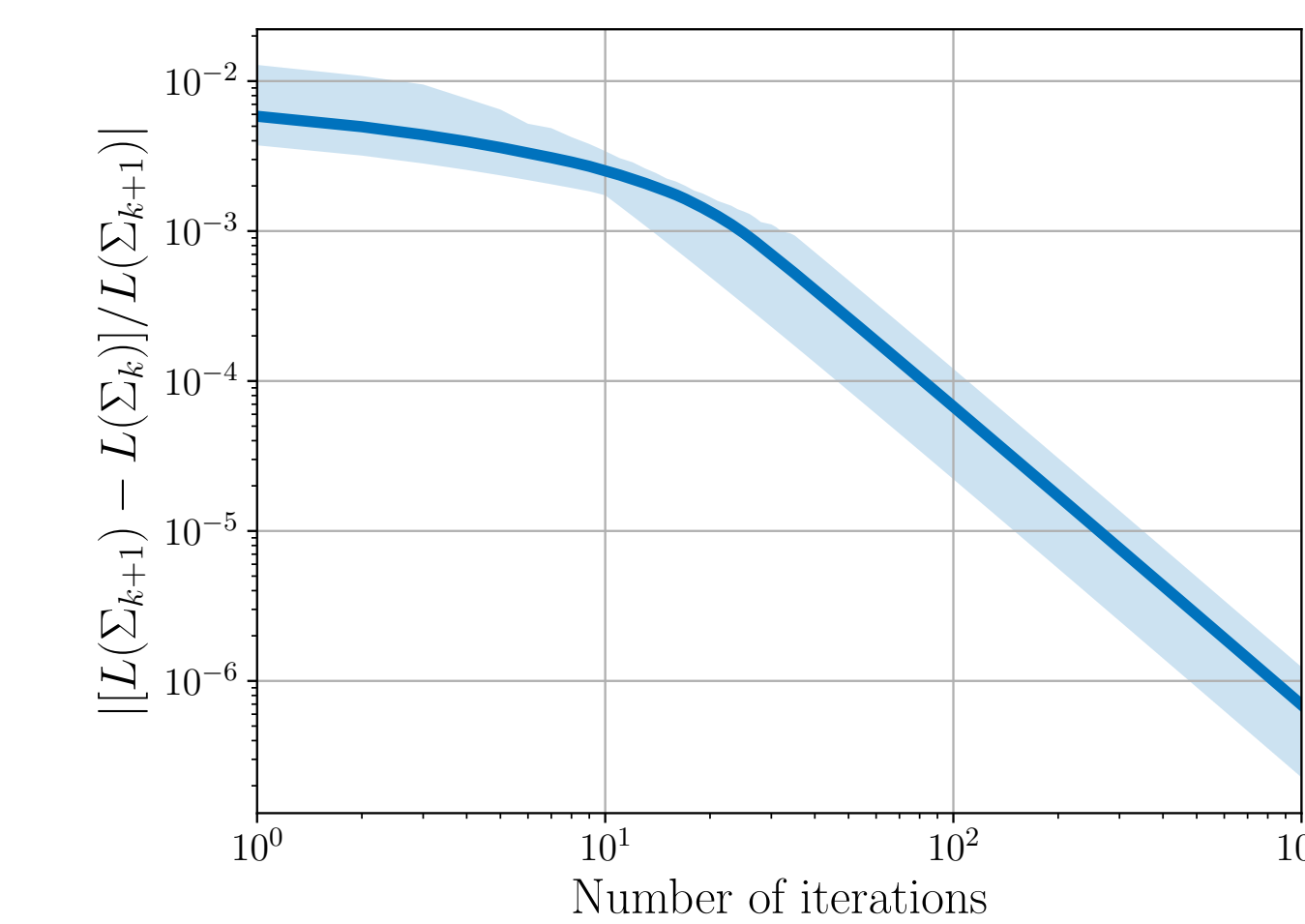
## Projected Geodesic Gradient Descent

1. Start from a feasible point $\Sigma$
2. Follow the Riemannian gradient $G$
3. Project back to $S_{++}^n$ manifold
$\Rightarrow$ use the exponential map $\text{Exp}_\Sigma(-\alpha G)$
4. Project back to $\mathcal{B}^{\text{FR}}$
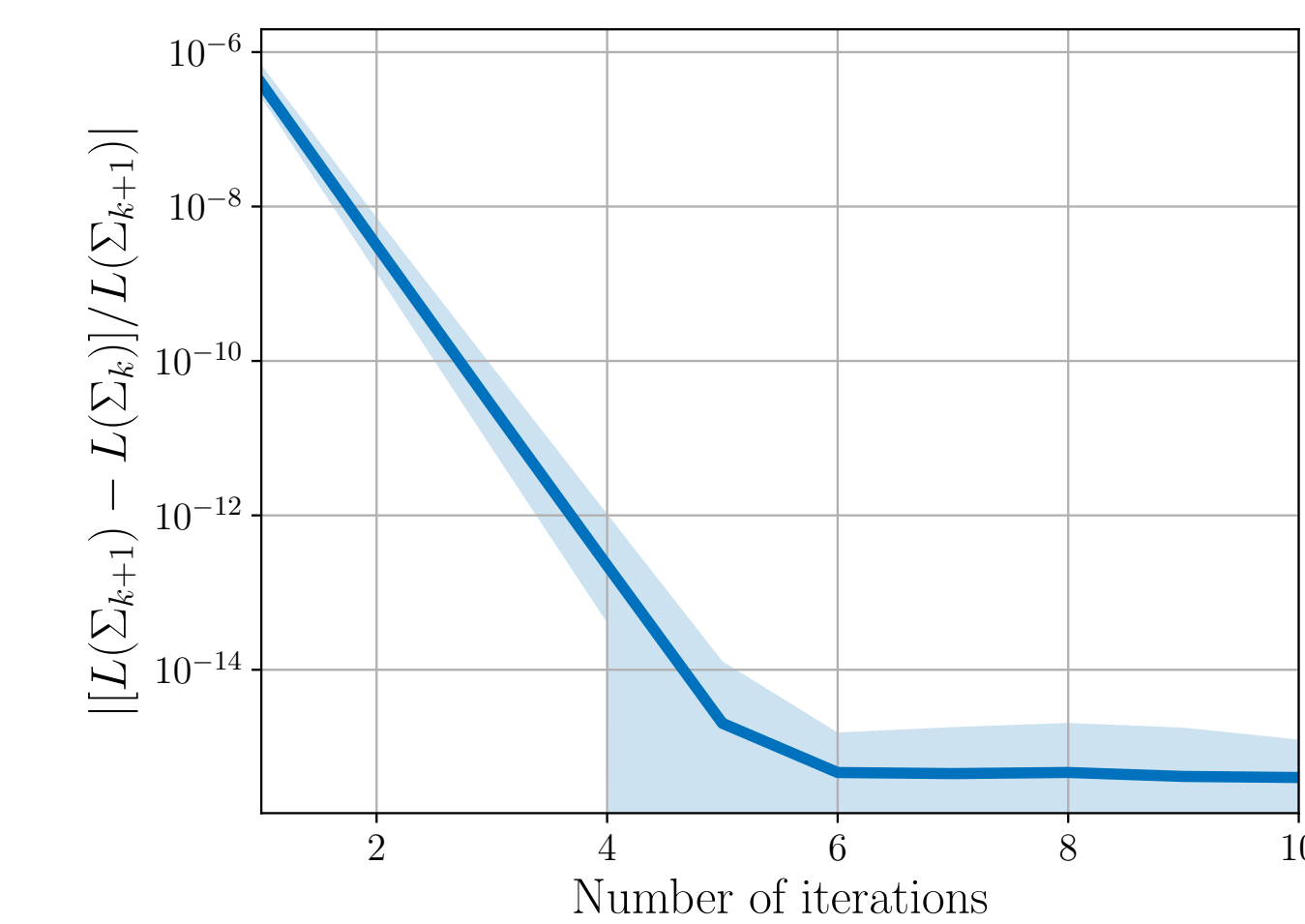


Based on [Zhang and Sra (2016)]

**Proposition 2.** If $d(\Sigma', \hat\Sigma) = \rho' > \rho$, then
$$\text{Proj}_{\mathcal{B}^{\text{FR}}}(\Sigma') = \hat\Sigma^{\frac{1}{2}}\left(\hat\Sigma^{-\frac{1}{2}}\Sigma'\hat\Sigma^{-\frac{1}{2}}\right)^{\frac{\rho}{\rho'}}\hat\Sigma^{\frac{1}{2}}$$

**Theorem 2.** With a constant stepsize $\alpha_k = \mathcal{O}(1/\sqrt{K})$, the projected geodesic gradient descent converges with the sublinear rate $\mathcal{O}(1/\sqrt{K})$.



(a) Convergence for $S \succeq 0$



(b) Convergence for $S \succ 0$

## OL Problem under the KL Divergence

The **OL** problem reduces to
$$\min_{\Sigma \succ 0} \text{Tr}\left[S\Sigma^{-1}\right] + \log\det \Sigma$$
$$\text{s.t. } \text{Tr}\left[\Sigma^{-1}\hat\Sigma\right] + \log\det(\Sigma\hat\Sigma^{-1}) - n \leq 2\rho.$$
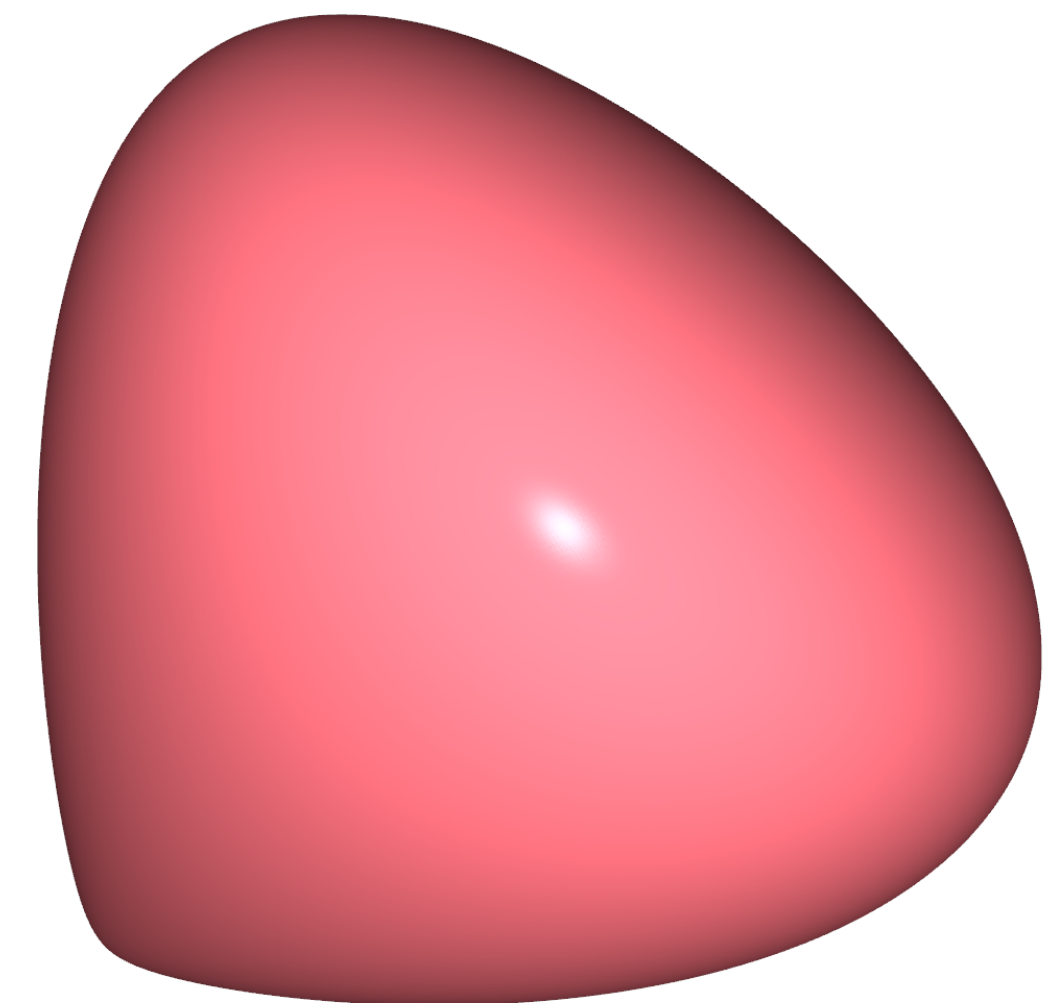
- Both objective and constraint are non-convex
- Convexification by substitution $X \leftarrow \Sigma^{-1}$
- Further reduction to one dimensional problem

**Theorem 3.** The **OL** problem is solved by
$$\Sigma^\star = S + \gamma^\star \hat\Sigma,$$
where $\gamma^\star$ is the solution of
$$\min_{\gamma^\star > 0} \gamma^\star(2\rho + \log\det \hat\Sigma) + n(1 + \gamma^\star)\log(1 + \gamma^\star) - (1 + \gamma^\star)\log\det(S + \gamma^\star \hat\Sigma).$$



## Flexible Discriminant Rules

- Classification problem with $Y \in \mathcal{C}, \mathcal{C} = \{1, \ldots, C\}$
- Bayes' Theorem implies that $\mathbb{P}(Y = c | X = x) \propto \pi_c \cdot f_c(x)$
- Assumption: $\hat{P}_c = \mathcal{N}(\hat\mu_c, \hat\Sigma_c)$ and $\hat\pi_c = N_c/N$
- $\hat\mu_c$ and $\hat\Sigma_c$ are estimated from training data
- QDA rule: $\mathcal{C}_{\text{QDA}}(x) \in \arg\max_{c \in \mathcal{C}}\left\{\frac{1}{2}\ell(x, \hat{P}_c) + \log(\hat\pi_c)\right\}$
- Our suggestion: $\mathcal{C}_{\text{flex}}(x) \in \arg\max_{c \in \mathcal{C}} \max_{P \in \mathcal{P}_c}\left\{\frac{1}{2}\ell(x, P) + \log(\hat\pi_c)\right\}$

**Empirical Experiments (UCI dataset):** average correct classification rates



## References

- C. Atkinson and A. F. Mitchell. Rao's distance measure. Sankhya: The Indian Journal of Statistics, Series A, 43(3):345-365, 1981.
- H. Zhang and S. Sra. First-order methods for geodesically convex optimization. In Conference on Learning Theory, pages 1617-1638, 2016.