# Negative curvature obstructs acceleration for g-convex optimization,
## even with exact first-order oracles

Chris Criscitiello

Nicolas Boumal

OPTIM, Chair of Continuous Optimization

Institute of Mathematics, EPFL

# Question

Is there a fully accelerated first-order algorithm for geodesically convex optimization with exact oracles?

# Question

Is there a fully accelerated first-order algorithm for geodesically convex optimization with exact oracles?

Short answer: No.

# Question

Is there a fully accelerated first-order algorithm for geodesically convex optimization with exact oracles?

Short answer: No.

Slightly longer answer: We show there are Riemannian manifolds and regimes where gradient descent is optimal (worst-case complexity).

# Question

Is there a fully accelerated first-order algorithm for geodesically convex optimization with exact oracles?

Short answer: No.

Slightly longer answer: We show there are Riemannian manifolds and regimes where gradient descent is optimal (worst-case complexity).

Builds on work of Hamilton and Moitra (2021), who show the answer is no when algorithms receive noisy information.

Hamilton and Moitra: "A No-Go Theorem for Acceleration in the Hyperbolic Plane" (2021)
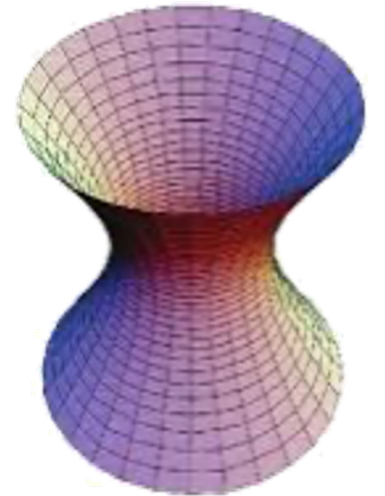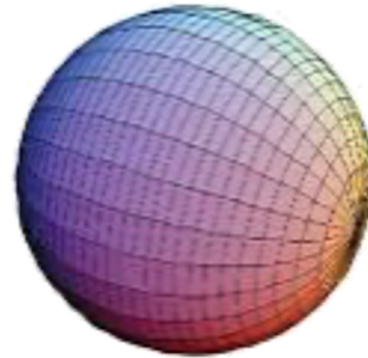
# Optimization on manifolds

$$\min_{x \in D \subset \mathcal{M}} f(x)$$

$\mathcal{M}$ is a Riemannian manifold

# Optimization on manifolds

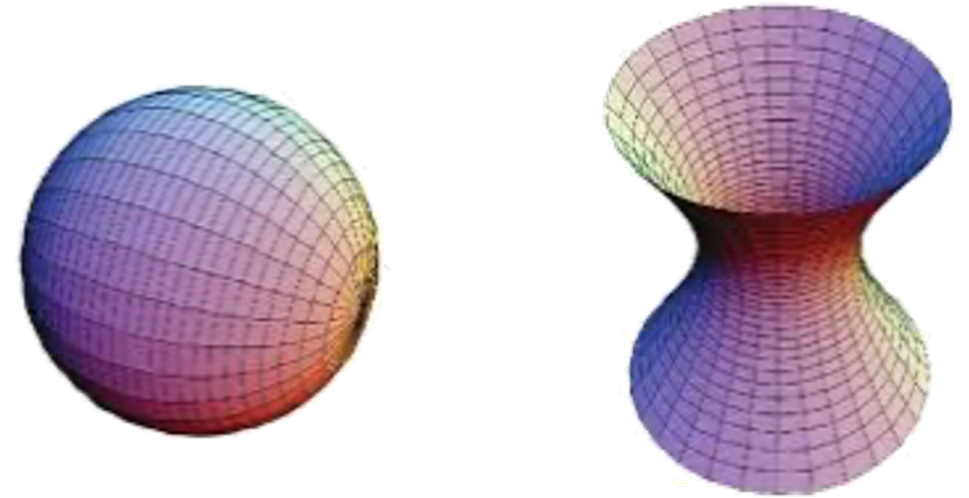$$\min_{x \in D \subset \mathcal{M}} f(x)$$
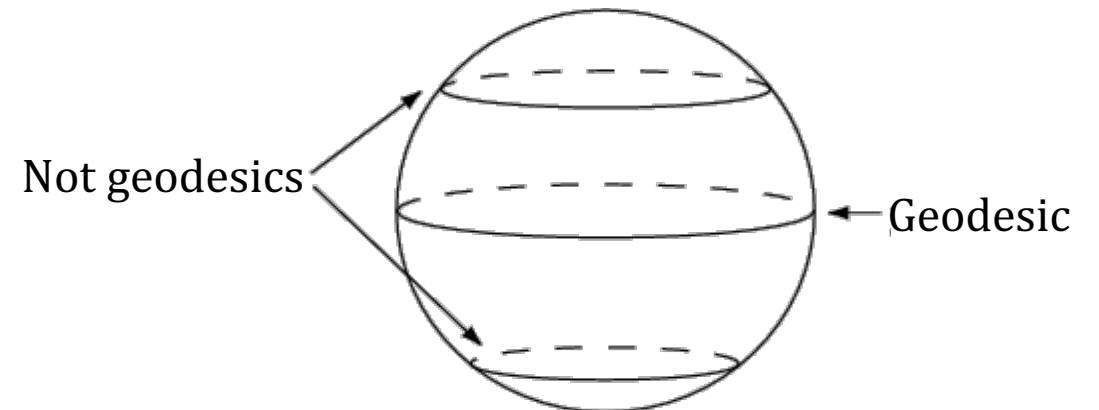
$\mathcal{M}$ is a Riemannian manifold

# Optimization on manifolds

$$\min_{x \in D \subset \mathcal{M}} f(x)$$

$\mathcal{M}$ is a Riemannian manifold

Geodesics: locally shortest paths.

Not geodesics

Geodesic

# Geodesically convex optimization

$$\min_{x \in D \subset \mathcal{M}} f(x)$$

Search space $D$ is a g-convex subset of a Riemannian manifold $\mathcal{M}$:

Cost $f$ is $\mu$-strongly g-convex:

# Geodesically convex optimization

$$\min_{x \in D \subset \mathcal{M}} f(x)$$

Search space $D$ is a g-convex subset of a Riemannian manifold $\mathcal{M}$:

For each $x, y \in D$, there is a unique minimizing geodesic $t \mapsto \gamma(t)$ contained in $D$, connecting $x, y$.
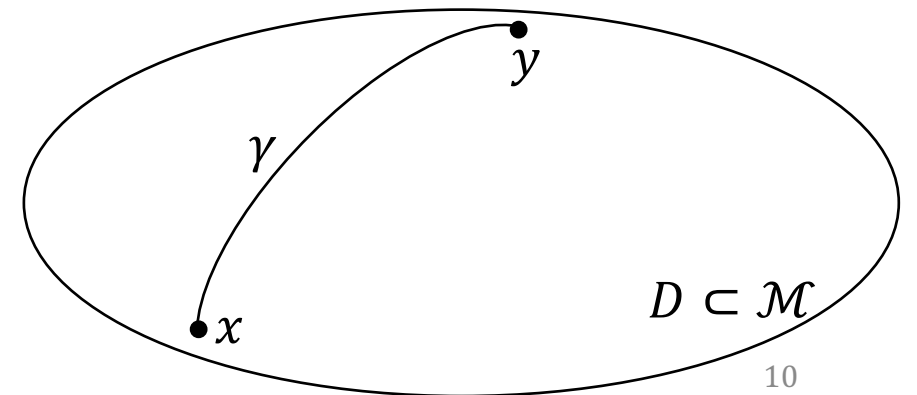
Cost $f$ is $\mu$-strongly g-convex:

# Geodesically convex optimization

$$\min_{x \in D \subset \mathcal{M}} f(x)$$

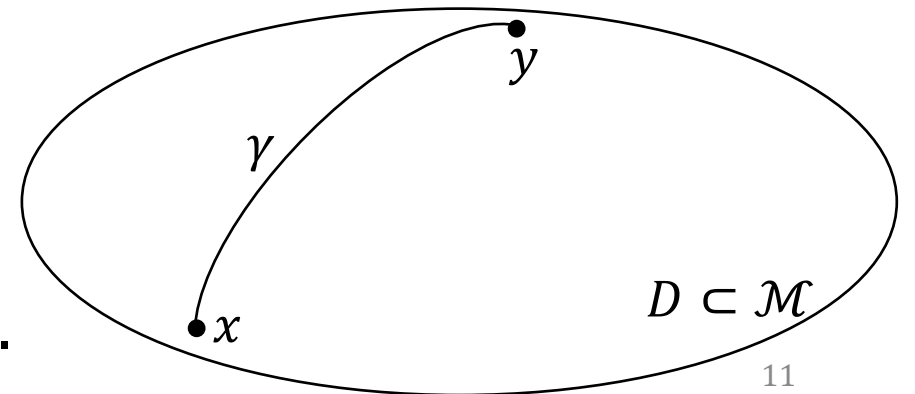Search space $D$ is a g-convex subset of a Riemannian manifold $\mathcal{M}$:

For each $x, y \in D$, there is a unique minimizing geodesic $t \mapsto \gamma(t)$ contained in $D$, connecting $x, y$.

Cost $f$ is $\mu$-strongly g-convex:

$$t \mapsto f(\gamma(t))$$

is $\mu$-strongly convex for any geodesic $\gamma$ in $D$.

# Hadamard manifolds

Complete, simply connected, with non-positive (intrinsic) curvature.

# Hadamard manifolds

Complete, simply connected, with non-positive (intrinsic) curvature.

Euclidean space: $\mathcal{M} = \mathbb{R}^d$

Hyperbolic space

# Hadamard manifolds

Complete, simply connected, with non-positive (intrinsic) curvature.

Euclidean space: $\mathcal{M} = \mathbb{R}^d$

Hyperbolic space

Positive definite matrices: $\mathcal{M} = \{P \in \mathbf{R}^{n \times n} : P = P^\top \text{ and } P \succ 0\}$
with affine-invariant metric $\langle X, Y \rangle_P = \text{Tr}(P^{-1} X P^{-1} Y)$.

# Hadamard manifolds

Complete, simply connected, with <span style="color:orange">non-positive (intrinsic) curvature</span>.

Euclidean space: $\mathcal{M} = \mathbb{R}^d$

Hyperbolic space

Positive definite matrices: $\mathcal{M} = \{P \in \mathbf{R}^{n \times n} : P = P^\top \text{ and } P \succ 0\}$
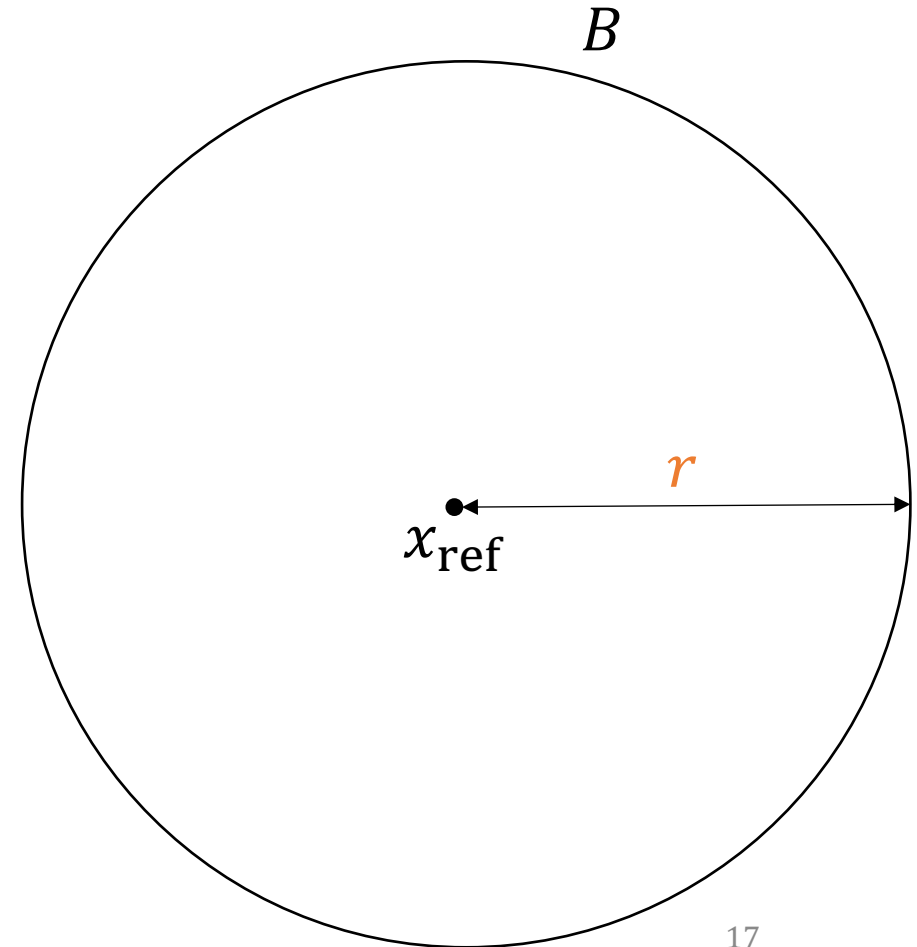with affine-invariant metric $\langle X, Y \rangle_P = \text{Tr}(P^{-1} X P^{-1} Y)$.

Non-example: Sphere

# Applications

- Fréchet mean (intrinsic averaging on Hadamard spaces) (e.g., Karcher)
- Gaussian mixture models (Hosseini + Sra)
- Optimistic likelihoods for Gaussians (Nguyen et al.)
- Robust Covariance estimation (Weisel + Zhang, Franks + Moitra)
- Metric learning (Zadeh et al.)
- Variants on PCA (Tang + Allen) [MLEs for matrix normal models]
- Operator/tensor scaling (Allen Zhu et al., Burgisser et al.)
  - Brascamp-Lieb constants, computational complexity, polynomial identity testing, hardness of robust subspace recovery, etc.
- Tree-like embeddings (Bacak)
- Sampling on Riemannian manifolds (Goyal + Shetty)

# Computational task

Geodesic ball $B = B(x_{\mathrm{ref}}, r)$ of radius $r$ in Hadamard space $\mathcal{M}$.

# Computational task

Geodesic ball $B = B(x_{\mathrm{ref}}, r)$ of radius $r$ in Hadamard space $\mathcal{M}$.

You know:
- $f$ is *L*-smooth in $B$ and $\mu$-strongly g-convex in $\mathcal{M}$;
- $f$ has a unique minimizer $x^*$ in $B$.

$B$

$x^*_\star$

$x_{\mathrm{ref}}$

$r$

Condition number $\kappa = \dfrac{L}{\mu}$

# Computational task

Geodesic ball $B = B(x_{\text{ref}}, r)$ of radius $r$ in Hadamard space $\mathcal{M}$.

You know:
- $f$ is $L$-smooth in $B$ and $\mu$-strongly g-convex in $\mathcal{M}$;
- $f$ has a unique minimizer $x^*$ in $B$.

You can query an oracle at $x$ to get $f(x), \text{grad } f(x)$
(exact info, no noise).



$B$

$x^*_\star$

$r$

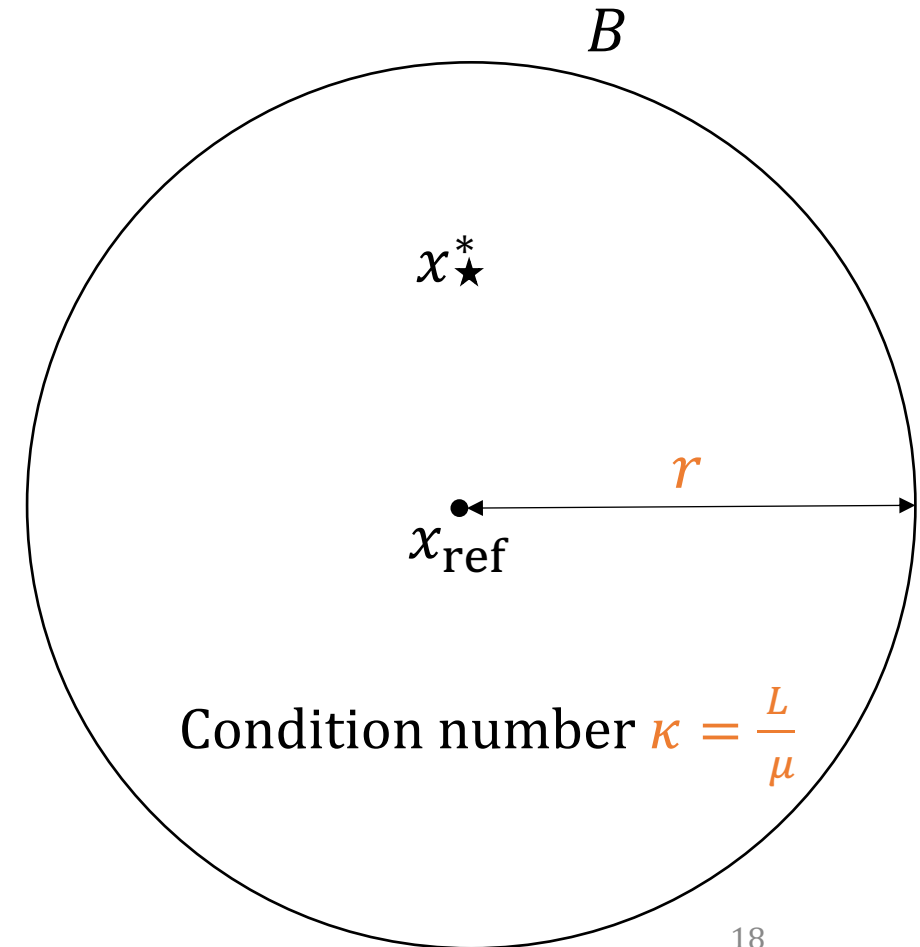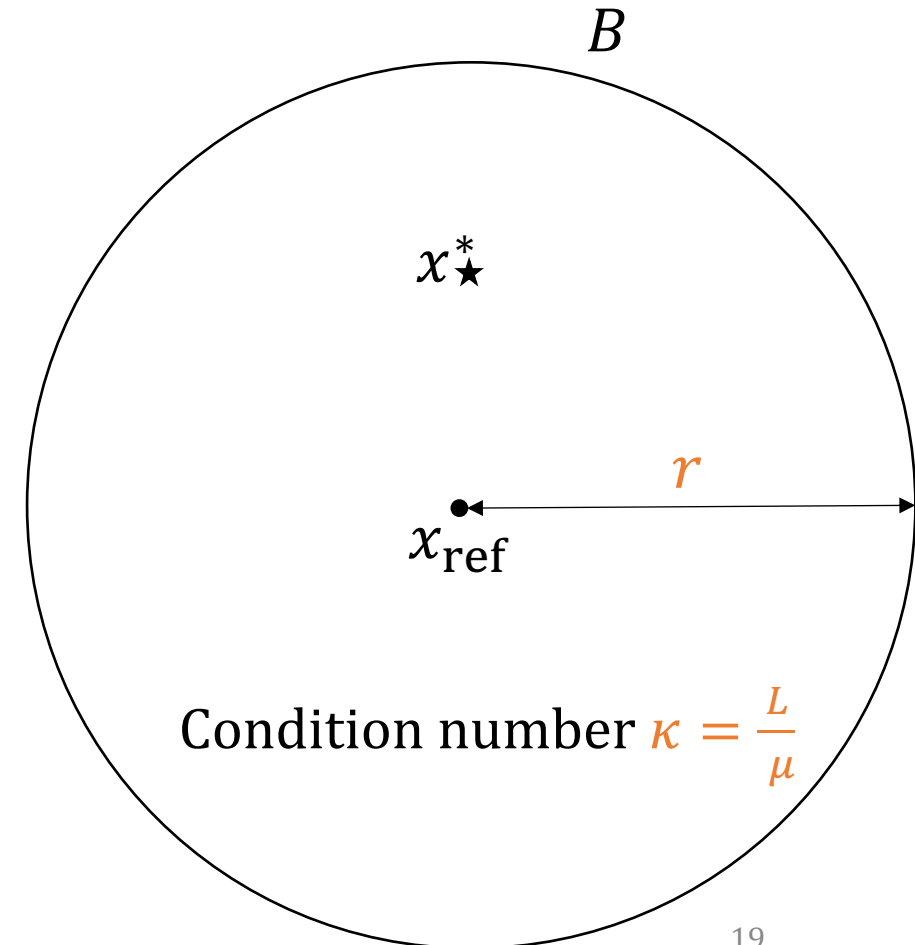$x_{\text{ref}}$

Condition number $\kappa = \dfrac{L}{\mu}$

# Computational task

Geodesic ball $B = B(x_{\text{ref}}, r)$ of radius $r$ in Hadamard space $\mathcal{M}$.

You know:
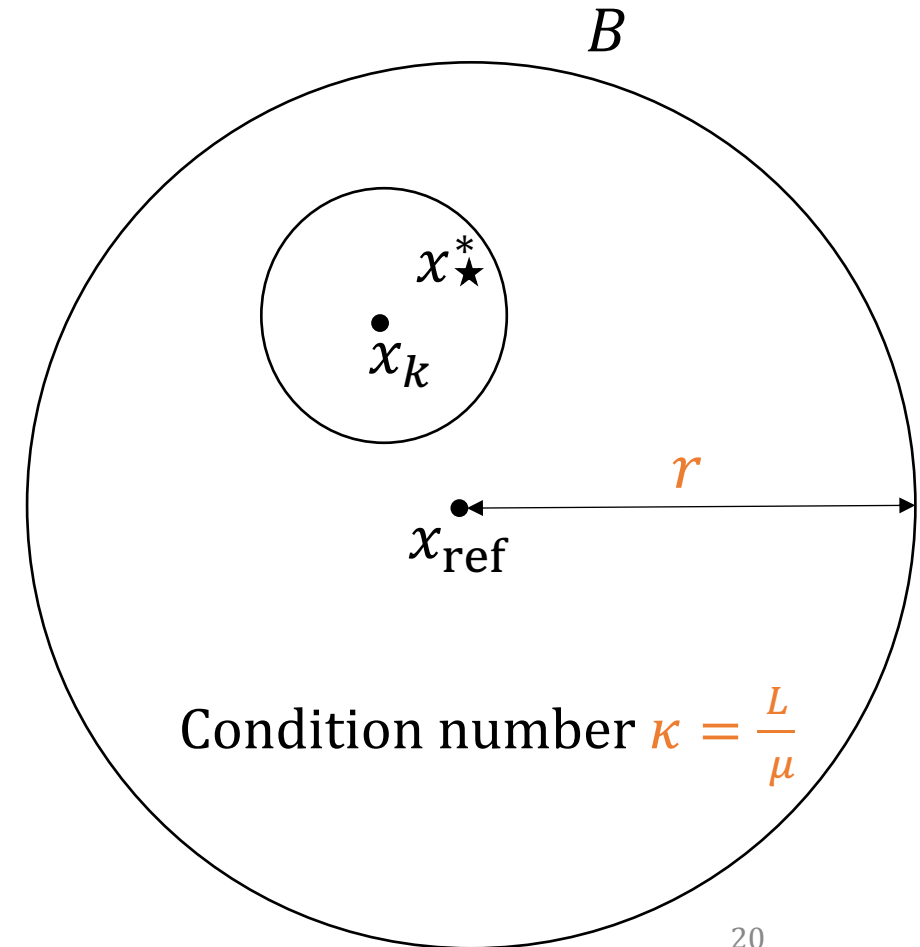- $f$ is $L$-smooth in $B$ and $\mu$-strongly g-convex in $\mathcal{M}$;
- $f$ has a unique minimizer $x^*$ in $B$.

You can query an oracle at $x$ to get $f(x), \operatorname{grad} f(x)$ (exact info, no noise).

Task: find a ball of radius $r/5$ containing $x^*$.

Least number of oracle queries necessary?



$B$

$x^*$

$x_k$

$r$

$x_{\text{ref}}$

Condition number $\kappa = \dfrac{L}{\mu}$

# What happens in $\mathbb{R}^d$?

If $\mathcal{M} = \mathbb{R}^d$:

Gradient Descent (GD)
$$x_{k+1} = x_k - \eta \nabla f(x_k)$$
$O(\kappa)$ oracle queries.



$B$

$x_\star^*$

$x_k$

$r$

$x_{\mathrm{ref}}$

Condition number $\kappa = \dfrac{L}{\mu}$

# What happens in $\mathbb{R}^d$?

If $\mathcal{M} = \mathbb{R}^d$:

Gradient Descent (GD)
$$x_{k+1} = x_k - \eta \nabla f(x_k)$$
$O(\kappa)$ oracle queries.

Nesterov's Accelerated Gradient method (NAG)
$$y_k = x_k + (1 - \theta)v_k$$
$$x_{k+1} = y_k - \eta \nabla f(y_k)$$
$$v_{k+1} = x_{k+1} - x_k$$
$\tilde{O}(\sqrt{\kappa})$ oracle queries.



$B$

$x^*$

$x_k$

$r$

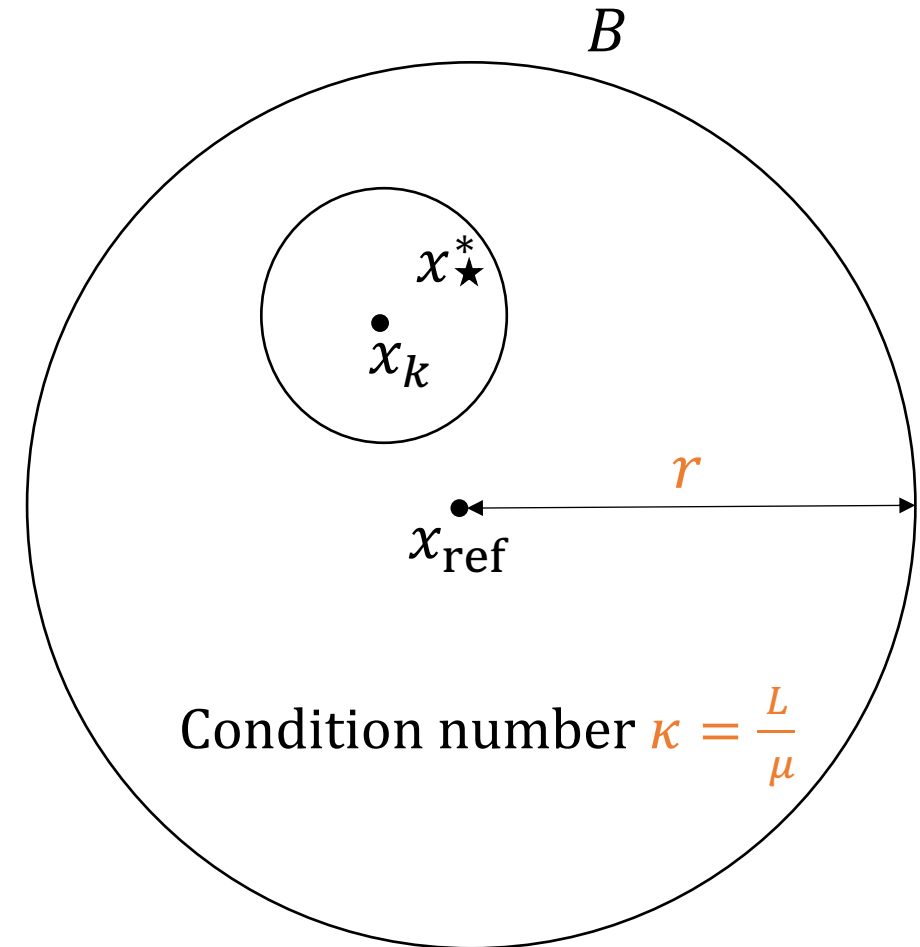$x_{\text{ref}}$

Condition number $\kappa = \dfrac{L}{\mu}$

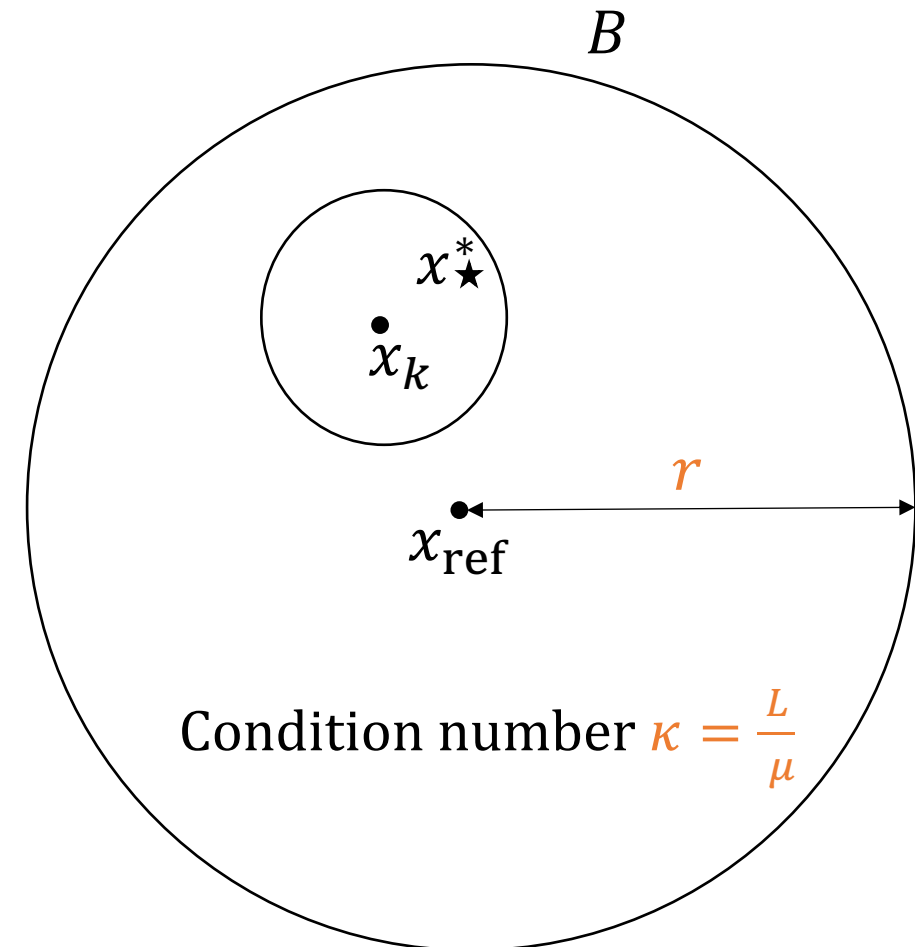# What happens in $\mathbb{R}^d$?

If $\mathcal{M} = \mathbb{R}^d$:

Gradient Descent (GD)
$$x_{k+1} = x_k - \eta \nabla f(x_k)$$
$O(\kappa)$ oracle queries.

Nesterov's Accelerated Gradient method (NAG)
$$y_k = x_k + (1-\theta)v_k$$
$$x_{k+1} = y_k - \eta \nabla f(y_k)$$
$$v_{k+1} = x_{k+1} - x_k$$
$\tilde{O}(\sqrt{\kappa})$ oracle queries.

NAG has optimal oracle complexity; GD does not.



$B$

$x^*$

$x_k$

$r$

$x_{\text{ref}}$

Condition number $\kappa = \dfrac{L}{\mu}$

# Optimal methods

What about on Riemannian manifolds?

# Optimal methods

What about on Riemannian manifolds?

Riemannian GD (RGD) requires $O(\kappa)$ oracle queries (when for example $\mathcal{M}$ is a hyperbolic space).

$$x_{k+1} = \exp_{x_k}(-\eta \operatorname{grad} f(x_k))$$

# Optimal methods

What about on Riemannian manifolds?

Riemannian GD (RGD) requires $O(\kappa)$ oracle queries (when for example $\mathcal{M}$ is a hyperbolic space).

$$x_{k+1} = \exp_{x_k}(-\eta \operatorname{grad} f(x_k))$$

Is there an algorithm using only $\tilde{O}(\sqrt{\kappa})$ queries in general (independent of $r$)?

# Optimal methods

What about on Riemannian manifolds?

Riemannian GD (RGD) requires $O(\kappa)$ oracle queries (when for example $\mathcal{M}$ is a hyperbolic space).

$$x_{k+1} = \exp_{x_k}(-\eta \operatorname{grad} f(x_k))$$

Is there an algorithm using only $\tilde{O}(\sqrt{\kappa})$ queries in general (independent of $r$)?

Partial positive result (Zhang, Ahn, Sra, Martinez-Rubio, Alimisis, et al.): you can accelerate in some cases (e.g., $r$ small).

# Main results

Let $\mathcal{M}$ be a Hadamard manifold of dimension $d \geq 2$ whose sectional curvatures are in the interval $[\mathrm{K}_{lo}, \mathrm{K}_{up}]$ with $\mathrm{K}_{up} < 0$.

Let $r = c_2 \, \kappa \, / \sqrt{-K_{lo}}$.

For hyperbolic spaces,
$K_{lo} = K_{up} = K < 0$

# Main results

Let $\mathcal{M}$ be a Hadamard manifold of dimension $d \geq 2$ whose sectional curvatures are in the interval $[\mathrm{K}_{lo}, \mathrm{K}_{up}]$ with $\mathrm{K}_{up} < 0$.

Let $r = c_2 \, \kappa \, / \sqrt{-K_{lo}}$.

For every deterministic algorithm $\mathcal{A}$, there is a $C^\infty$ function $f$ which is

- 1-strongly g-convex in all of $\mathcal{M}$;
- $\kappa$-smooth in the geodesic ball $B(x_{\mathrm{ref}}, r)$;
- and has (unique) minimizer in $B(x_{\mathrm{ref}}, r)$;

# Main results

Let $\mathcal{M}$ be a Hadamard manifold of dimension $d \geq 2$ whose sectional curvatures are in the interval $[\mathrm{K}_{lo}, \mathrm{K}_{up}]$ with $\mathrm{K}_{up} < 0$.
Let $r = c_2\,\kappa\,/\sqrt{-K_{lo}}$.

For every deterministic algorithm $\mathcal{A}$, there is a $C^\infty$ function $f$ which is
- 1-strongly g-convex in all of $\mathcal{M}$;
- $\kappa$-smooth in the geodesic ball $B(x_{\mathrm{ref}}, r)$;
- and has (unique) minimizer in $B(x_{\mathrm{ref}}, r)$;

such that algorithm $\mathcal{A}$ requires at least

$$\Omega\left(\sqrt{\frac{K_{up}}{K_{lo}}}\,\frac{\kappa}{\log \kappa}\right)$$

queries in order to find a point $x \in \mathcal{M}$ within $r/5$ of the minimizer of $f$.

# Main results

Let $\mathcal{M}$ be a Hadamard manifold of dimension $d \geq 2$ whose sectional curvatures are in the interval $[K_{lo}, K_{up}]$ with $K_{up} < 0$.

Let $r = c_2 \kappa / \sqrt{-K_{lo}}$.

For every deterministic algorithm $\mathcal{A}$, there is a $C^\infty$ function $f$ which is

- 1-strongly g-convex in all of $\mathcal{M}$;
- $\kappa$-smooth in the geodesic ball $B(x_{\mathrm{ref}}, r)$;
- and has (unique) minimizer in $B(x_{\mathrm{ref}}, r)$;

such that algorithm $\mathcal{A}$ requires at least

$$\Omega\left(\sqrt{\frac{K_{up}}{K_{lo}}} \frac{\kappa}{\log \kappa}\right) \implies \begin{array}{l} O(\sqrt{\kappa}) \text{ rate is impossible;} \\ \text{RGD is optimal (up to log).} \end{array}$$

queries in order to find a point $x \in \mathcal{M}$ within $r/5$ of the minimizer of $f$.

# Other settings

$n \times n$ positive definite matrices with affine-invariant metric.

Smooth nonstrongly g-convex optimization ($\mu = 0$).
        There are regimes where GD is optimal.

Nonsmooth g-convex optimization.

# Proof idea

Geodesic ball volumes grow much faster in negatively curved spaces than flat spaces.

# Proof idea

Geodesic ball volumes grow much faster in negatively curved spaces than flat spaces.

It's harder to find a point in a ball just because there's so much more space to explore.

* First highlighted by Hamilton and Moitra.

# Proof idea

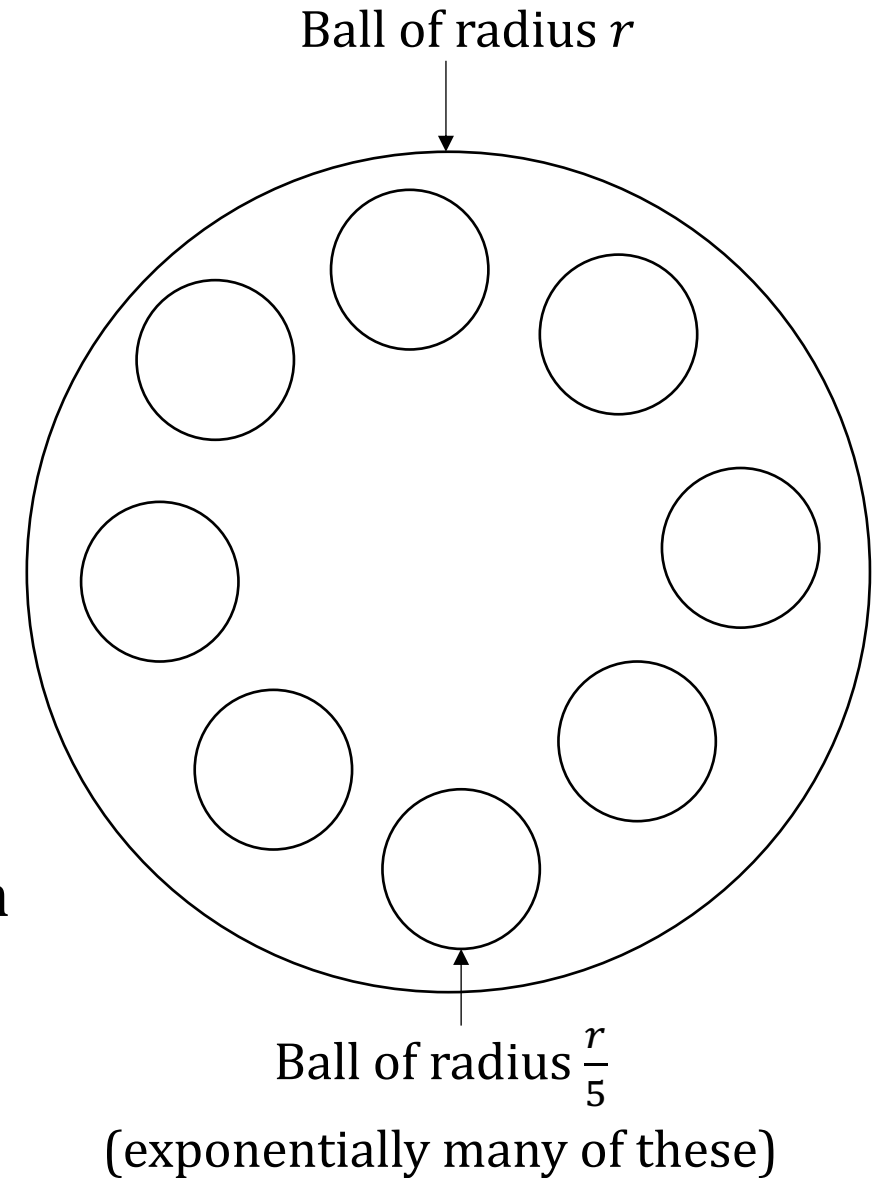Geodesic ball volumes grow much faster in negatively curved spaces than flat spaces.

It's harder to find a point in a ball just because there's so much more space to explore.

\* First highlighted by Hamilton and Moitra.

How many disjoint balls of radius $r/5$ contained in every ball of radius $r$?
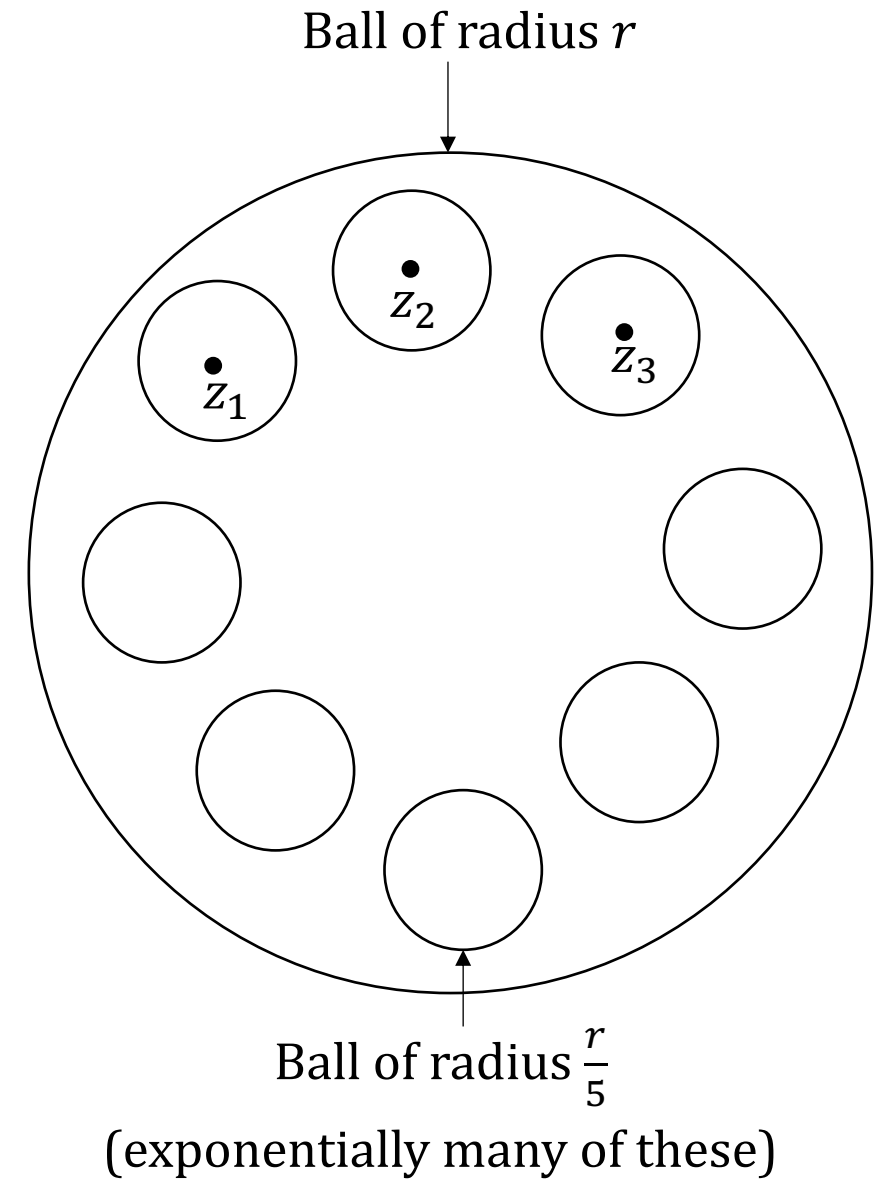
$e^{\Theta(rd)}$ in hyperbolic space

$e^{\Theta(d)}$ in Euclidean space

Ball of radius $r$

Ball of radius $\frac{r}{5}$

(exponentially many of these)

# Proof idea

Start with

$$x \mapsto \frac{1}{2}\text{dist}\left(x, z_j\right)^2, j = 1, \dots, N$$

Ball of radius $r$

$z_1$

$z_2$

$z_3$

Ball of radius $\frac{r}{5}$

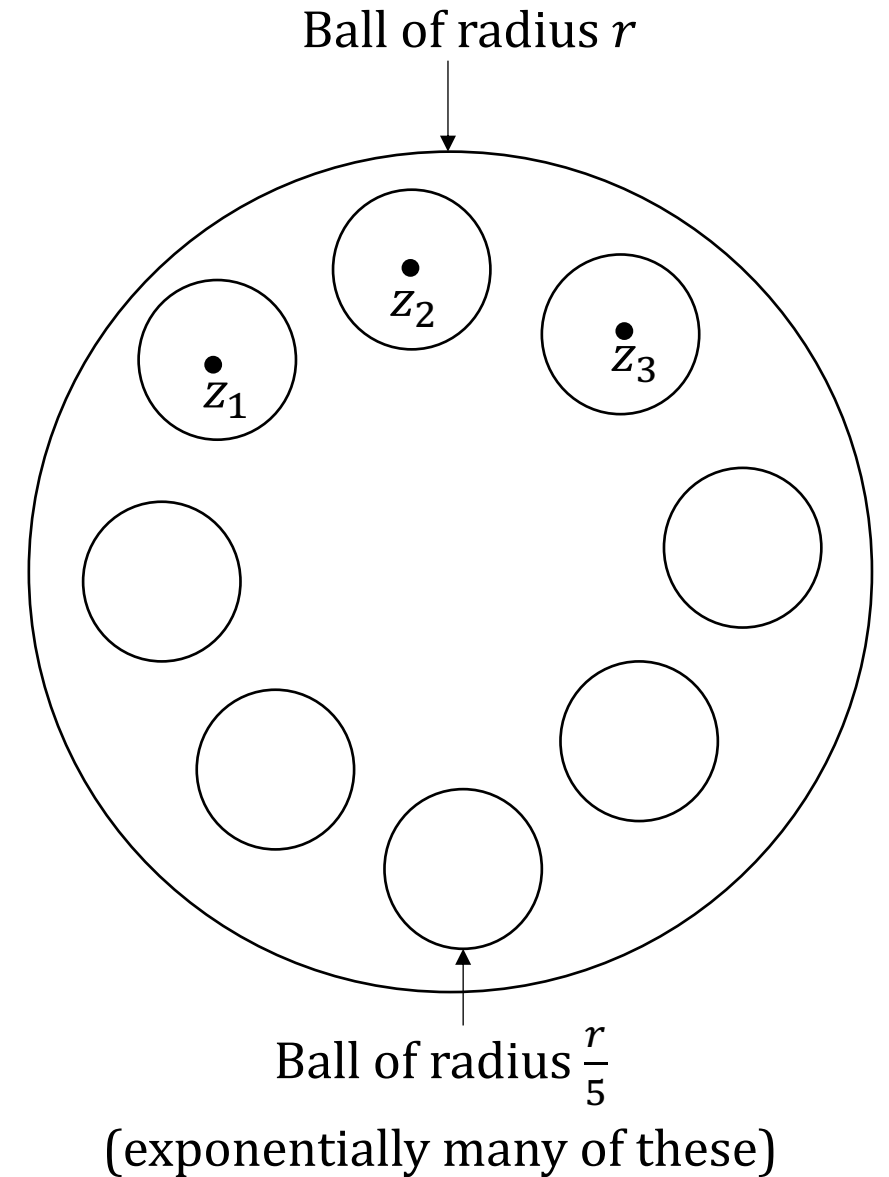(exponentially many of these)

# Proof idea

Start with

$$x \mapsto \frac{1}{2}\operatorname{dist}(x, z_j)^2, j = 1, \ldots, N$$

Gradients of these functions point directly towards the minimizer.



Ball of radius $r$

Ball of radius $\frac{r}{5}$

(exponentially many of these)

# Proof idea

Our solution: Add perturbations

$$x \mapsto \frac{1}{2}\text{dist}(x, z_j)^2 + H_{j,k}(x), \qquad \left\|\text{Hess}\, H_{j,k}(x)\right\| \leq \frac{1}{2}.$$

# Proof idea

Our solution: Add perturbations

$$x \mapsto \frac{1}{2}\text{dist}(x, z_j)^2 + H_{j,k}(x), \qquad \left\| \text{Hess } H_{j,k}(x) \right\| \leq \frac{1}{2}.$$

The perturbation $H_{j,k}$ is constructed adversarially using a resisting oracle.
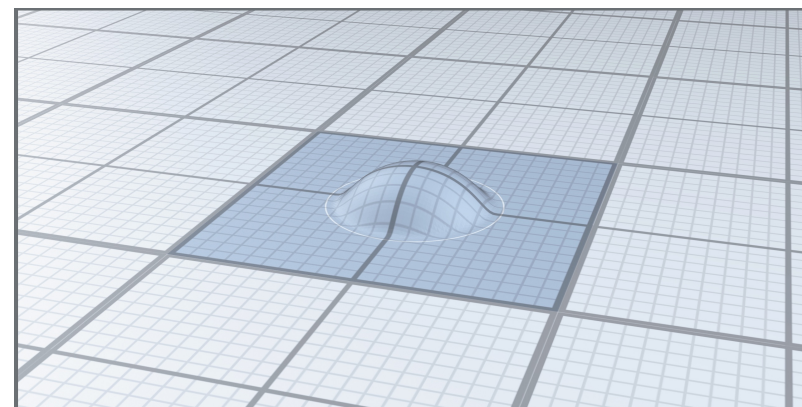
# Proof idea

Our solution: Add perturbations

$$x \mapsto \frac{1}{2}\operatorname{dist}(x, z_j)^2 + H_{j,k}(x), \qquad \left\|\operatorname{Hess} H_{j,k}(x)\right\| \leq \frac{1}{2}.$$

The perturbation $H_{j,k}$ is constructed adversarially using a resisting oracle.

Perturbation is a sum of bump functions

$$H_{j,k}(x) = \sum_{m=1}^{k} h_{j,m}.$$

# Future directions

Tighter upper/lower bounds, e.g., Kim and Yang (2022)

"Accelerated Gradient Methods for Geodesically Convex Optimization: Tractable Algorithms and Convergence Analysis"
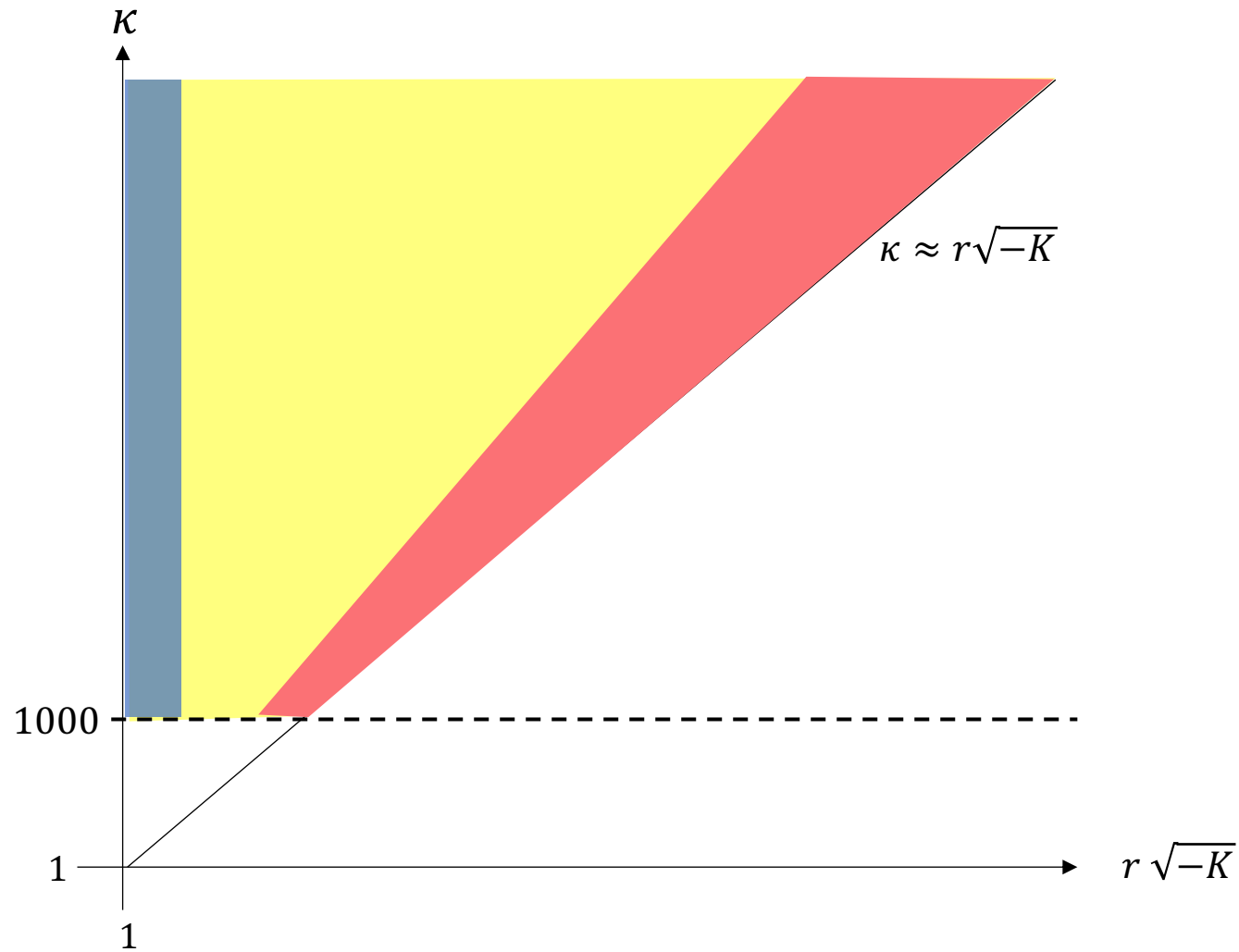
Randomized algorithms which receive exact information?

Ellipsoid method?

Interior-point methods?

# Appendix

# What we know (for hyperbolic spaces)



$\kappa \approx r\sqrt{-K}$

# Main results

$n \times n$ positive definite matrices with affine-invariant metric

# Main results

$n \times n$ positive definite matrices with affine-invariant metric

It is Hadamard, but does not satisfy assumptions of previous theorem: sectional curvature can be zero.

# Main results

$n \times n$ positive definite matrices with affine-invariant metric

It is Hadamard, but does not satisfy assumptions of previous theorem: sectional curvature can be zero.

Still, can prove the lower bound $\Omega\left(\frac{1}{n}\frac{\kappa}{\log\kappa}\right)$.

# Nonstrongly g-convex case

Can also show that acceleration is impossible in the nonstrongly g-convex case ($\mu = 0$).

# Nonstrongly g-convex case

Can also show that acceleration is impossible in the nonstrongly g-convex case ($\mu = 0$).

Have the lower bound $\Omega\left(\frac{1}{\epsilon} \cdot \frac{1}{\log^3(\epsilon^{-1})}\right)$ for finding a point $x$ with $f(x) - f(x^*) \leq \epsilon$.

Means a version of RGD is optimal.

# Nonstrongly g-convex case

Can also show that acceleration is impossible in the nonstrongly g-convex case ($\mu = 0$).

Have the lower bound $\Omega\left(\frac{1}{\epsilon} \cdot \frac{1}{\log^3(\epsilon^{-1})}\right)$ for finding a point $x$ with $f(x) - f(x^*) \leq \epsilon$.

Means a version of RGD is optimal.

Compare with NAG, which uses at most $O\left(\frac{1}{\sqrt{\epsilon}}\right)$ queries in Euclidean spaces.

# Application: robust covariance estimation

IID samples $x_i \in \mathbb{R}^p, i = 1, \ldots, n$, coming from an elliptical distribution:
$$x \sim u \, \Sigma^{1/2} v$$

where $\Sigma > 0$ is fixed (the shape matrix), $u$ is a scalar r.v., and $v \sim \mathbb{S}^{p-1}$.

# Application: robust covariance estimation

IID samples $x_i \in \mathbb{R}^p, i = 1, \ldots, n,$ coming from an elliptical distribution:
$$x \sim u\,\Sigma^{1/2} v$$

where $\Sigma > 0$ is fixed (the shape matrix), $u$ is a scalar r.v., and $v \sim \mathbb{S}^{p-1}$.

Tyler's M-estimator for the shape matrix:
$$\hat{\Sigma} = \operatorname*{argmin}_{\Sigma > 0,\ \mathrm{Tr}(\Sigma) = p} \frac{p}{n} \sum_{i=1}^{n} \log\left(x_i^\top \Sigma^{-1} x_i\right) + \log\det(\Sigma)$$

Can also be derived as an MLE.

# Application: robust covariance estimation

IID samples $x_i \in \mathbb{R}^p, i = 1, \ldots, n$, coming from an elliptical distribution:
$$x \sim u \, \Sigma^{1/2} v$$

where $\Sigma \succ 0$ is fixed (the shape matrix), $u$ is a scalar r.v., and $v \sim \mathbb{S}^{p-1}$.

Tyler's M-estimator for the shape matrix:
$$\hat{\Sigma} = \underset{\Sigma \succ 0, \; \mathrm{Tr}(\Sigma)=p}{\mathrm{argmin}} \frac{p}{n} \sum_{i=1}^{n} \log\left(x_i^\top \Sigma^{-1} x_i\right) + \log \det(\Sigma)$$

Is g-convex for PD matrices (with affine-invariant metric).

$\rightarrow$ new algorithms/analysis + analysis for Tyler's iterative procedure

# Application: robust covariance estimation

IID samples $x_i \in \mathbb{R}^p, i = 1, \ldots, n$, coming from an elliptical distribution:

$$x \sim u \, \Sigma^{1/2} v$$

where $\Sigma > 0$ is fixed (the shape matrix), $u$ is a scalar r.v., and $v \sim \mathbb{S}^{p-1}$.

Tyler's M-estimator for the shape matrix:

$$\hat{\Sigma} = \underset{\Sigma > 0, \ \mathrm{Tr}(\Sigma) = p}{\mathrm{argmin}} \frac{p}{n} \sum_{i=1}^{n} \log\left(x_i^\top \Sigma^{-1} x_i\right) + \log \det(\Sigma)$$

Is a specific instance of the operator scaling problem.

Sources: Weisel + Zhang, Franks + Moitra