

Overview

- We recover the **asymptotic equivalence** of the Rao, Wald, and likelihood-ratio tests from a **non-asymptotic viewpoint**.
- We characterize the **critical sample size** beyond which the equivalence holds asymptotically under the **null** hypotheses.
- We also analyze the **statistical power** under both the **fixed and alternative** hypotheses.
- We establish an estimation bound that matches the **misspecified Cramér-Rao** lower bound.

Goodness-of-Fit Testing

Problem. Let $Z \sim \mathbb{P}$ and $\mathcal{P}_\Theta := \{P_\theta : \theta \in \Theta \subset \mathbb{R}^d\}$. Assume there exists $\theta_* \in \Theta$ such that $\mathbb{P} = P_{\theta_*}$. Given an i.i.d. sample $\{Z_i\}_{i=1}^n$, we want to infer properties of θ_* via

$$\mathcal{H}_0 : \theta_* = \theta_0 \leftrightarrow \mathcal{H}_1 : \theta_* \neq \theta_0.$$

- A **test statistic** $T := T(Z_1, \dots, Z_n)$ and a **prescribed critical value** t_n .
 - **Reject** the null \mathcal{H}_0 if $T > t_n$.
 - **Type I error rate** $\Pr(T > t_n | \mathcal{H}_0)$ and **statistical power** $\Pr(T > t_n | \mathcal{H}_1)$.
- Notation.** Loss function $\ell(\theta; z) := -\log P_\theta(z)$.
- **Population risk** $L(\theta) := \mathbb{E}[-\log P_\theta(Z)]$ and **empirical risk** $L_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; Z_i)$.
 - **Empirical risk minimizer** $\hat{\theta}_n := \arg \min_{\theta \in \Theta} L_n(\theta)$.
 - **Gradient** $S(\theta; z) := \nabla_\theta \ell(\theta; z)$, $S(\theta) := \mathbb{E}[S(\theta; Z)] = \nabla_\theta L(\theta)$, and $S_n(\theta) := \frac{1}{n} \sum_{i=1}^n S(\theta; Z_i)$.
 - **Hessian** $H(\theta; z) := \nabla_\theta^2 \ell(\theta; z)$, $H(\theta) := \mathbb{E}[H(\theta; Z)]$, and $H_n(\theta) := \frac{1}{n} \sum_{i=1}^n H(\theta; Z_i)$.

Three classical goodness-of-fit tests.

- **Rao test**— $T_{\text{Rao}} := S_n(\hat{\theta}_n)^\top H_n(\hat{\theta}_n)^{-1} S_n(\hat{\theta}_n)$.
- **Wald test**— $T_{\text{Wald}} := (\hat{\theta}_n - \theta_0)^\top H_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0)$.
- **The likelihood-ratio test**— $T_{\text{LR}} := 2[L_n(\hat{\theta}_n) - L_n(\theta_0)]$.

Asymptotic equivalence of the three tests.

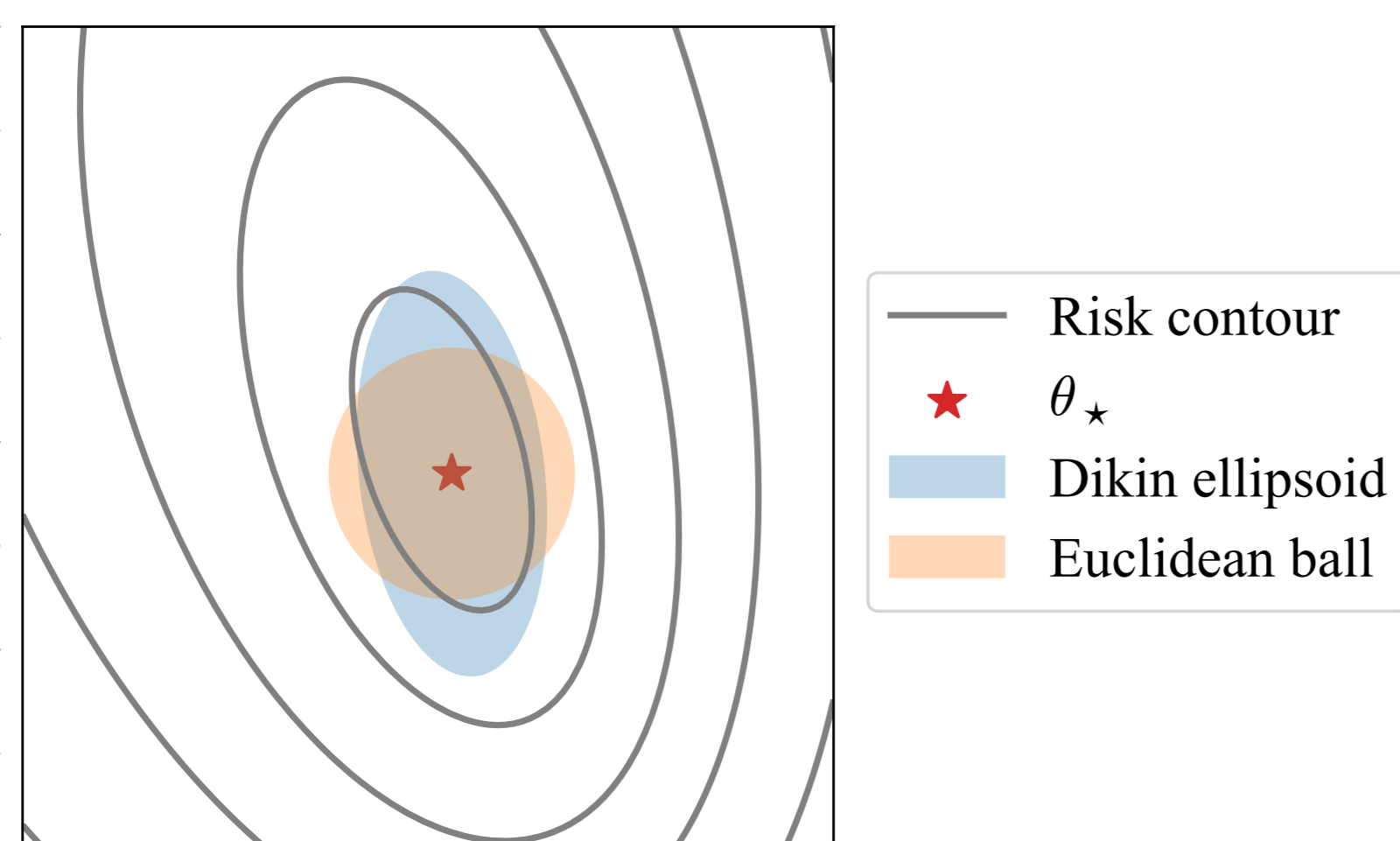
- Asymptotic distribution.
 - $\sqrt{n}S_n(\hat{\theta}_n) \rightarrow_d \mathcal{N}_d(0, G(\theta_0))$ under \mathcal{H}_0 .
 - Well-specified model, i.e., $\mathbb{P} \in \mathcal{P}_\Theta$, implies that $G(\theta_0) = H(\theta_0)$.
 - $nT_{\text{Rao}} \rightarrow_d \chi_d^2$ under \mathcal{H}_0 .
- Asymptotic equivalence.
 - $S_n(\hat{\theta}_n) = S_n(\theta_0) - S_n(\hat{\theta}_n) = H_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0)$.
 - $T_{\text{Rao}} = (\hat{\theta}_n - \theta_0)^\top H_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) + o_p(1) = T_{\text{Wald}} + o_p(1)$.
 - $T_{\text{LR}} = 2S_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) + (\hat{\theta}_n - \theta_0)^\top H_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) = T_{\text{Wald}} + o_p(1)$.

Preliminaries

Dikin ellipsoid. A Dikin ellipsoid at θ_* of radius r

$$\Theta_r(\theta_*) := \{\theta \in \Theta : \|H_*^{1/2}(\theta - \theta_*)\| < r\}.$$

- The shape of a **Euclidean ball** is **always the same**.
- The shape of a **Dikin ellipsoid** is **adapted to the geometry** near the optimum.



Generalized self-concordance. Let f be convex, $R > 0$, and $\nu > 0$. We say f is (R, ν) -generalized self-concordant if (on a high level)

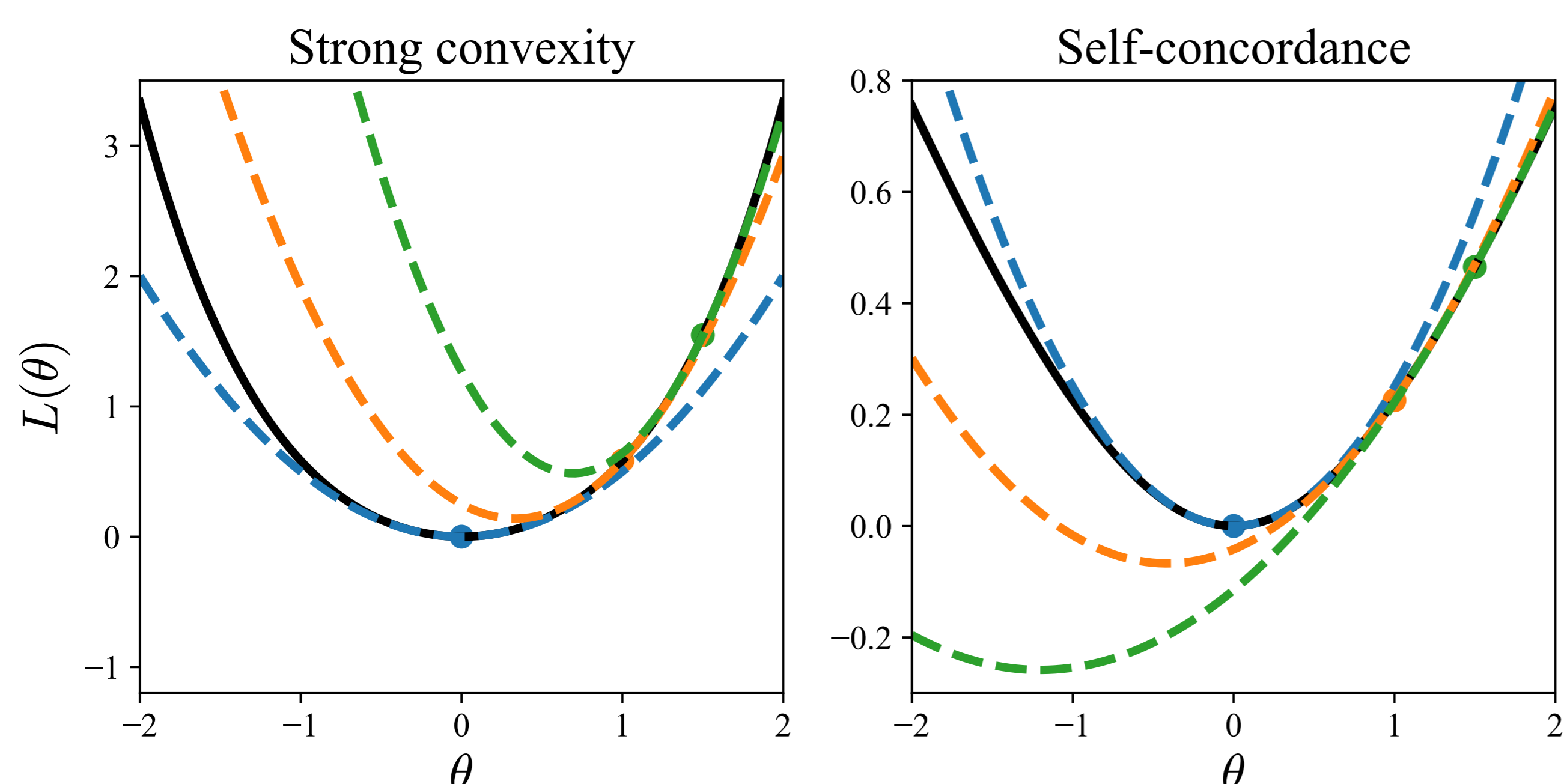
$$\|\nabla^3 f(x)\| \lesssim R \|\nabla^2 f(x)\|^\nu.$$

We give two examples of losses as functions of parameters.

- For **linear regression**, its loss is (R, ν) -generalized self-concordant for any $R > 0$ and $\nu > 0$.
- For **logistic regression** with $\|X\|_2 \leq_{a.s.} M$, its loss is $(2M, 2)$ -generalized self-concordant.

Strong convexity Self-concordance

Hessian lower bound	Global	Local
Hessian varying rate	No control	Slow



Concentration of Hessian. A key result towards deriving our bounds is

$$[1 - c_n(\delta)]H(\theta) \preceq H_n(\theta) \preceq [1 + c_n(\delta)]H(\theta)$$

with probability at least $1 - \delta$, where $c_n(\delta) = O(\sqrt{\log(d/\delta)/n})$.

Main Results

Type I error rate. Let $\theta_* = \theta_0$. We have, with probability at least $1 - \delta$,

$$nT_{\text{Rao}} \lesssim d + \log \frac{e}{\delta}, \quad \text{whenever } n \gtrsim \log \frac{2d}{\delta}.$$

Additionally, with $\lambda_* := \lambda_{\min}(H_*)$,

$$nT_{\text{Wald}}, nT_{\text{LR}} \lesssim d + \log \frac{e}{\delta}, \quad \text{whenever } n \gtrsim \log \frac{2d}{\delta} + d \frac{R^2}{\lambda_*^{3-\nu}}.$$

- Demonstrate that the three tests have a tail behavior that is governed by a χ_d^2 distribution.
- Characterize the **critical sample size** enough to enter the asymptotic regime.

Statistical power. Let $\theta_* \rightarrow_{n \rightarrow \infty} \theta_0$. Let $t_n(\alpha)$ be the $(1 - \alpha)$ -quantile of χ_d^2 . Let $\Omega(\theta) := G(\theta)^{1/2} H(\theta)^{-1} G(\theta)^{1/2}$ and $h(\tau) := \min\{\tau^2, \tau\}$. The following statements hold for sufficiently large n .

- Let $\tau_n \asymp \|H(\theta_0)^{1/2}(\theta_* - \theta_0)\|^2$. We have

$$\Pr(T_{\text{Rao}} > t_n(\alpha)) \begin{cases} \leq 2de^{-Cn} + e^{-Ch(\|\Omega(\theta_0)\|_2^{-1})} & \text{if } \|H(\theta_0)^{1/2}(\theta_* - \theta_0)\| = O(n^{-1/2}) \\ \geq 1 - 2de^{-Cn} - e^{-Ch(n\tau_n\|\Omega(\theta_0)\|_2^{-1})} & \text{if } \|H(\theta_0)^{1/2}(\theta_* - \theta_0)\| = \omega(n^{-1/2}). \end{cases}$$

- Let $\tau'_n \asymp \|H(\theta_0)^{1/2}(\theta_* - \theta_0)\|^2$. We have

$$\Pr(T_{\text{Wald}}, T_{\text{LR}} > t_n(\alpha)) \begin{cases} \leq 2nde^{-C(\frac{\lambda_*^{\nu-1}}{R^2 d})} + e^{-Ch(\|\Omega(\theta_0)\|_2^{-1})} & \text{if } \|H(\theta_0)^{1/2}(\theta_* - \theta_0)\| = O(n^{-1/2}) \\ \geq 1 - 2nde^{-C(\frac{\lambda_*^{\nu-1}}{R^2 d})} - e^{-Ch(n\tau'_n\|\Omega(\theta_0)\|_2^{-1})} & \text{if } \|H(\theta_0)^{1/2}(\theta_* - \theta_0)\| = \omega(n^{-1/2}). \end{cases}$$

To summarize,

- If $\|H(\theta_0)^{1/2}(\theta_* - \theta_0)\| = O(n^{-1/2})$, the power is asymptotically **upper bounded by a constant**.
- If $\|H(\theta_0)^{1/2}(\theta_* - \theta_0)\| = \omega(n^{-1/2})$, the power **tends to one** at rate $O(\exp(-n\|H(\theta_0)^{1/2}(\theta_* - \theta_0)\|^2))$.

Estimation Bound under Model Misspecification

Estimation bound. Assume $\mathbb{P} \notin \mathcal{P}_\Theta$ and let $\theta_* := \arg \min_{\theta \in \Theta} L(\theta)$. It holds that

$$\|H_n(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_*)\| \lesssim \frac{d_*}{n} + \frac{\|\Omega(\theta_*)\|_2}{n} \log \frac{e}{\delta},$$

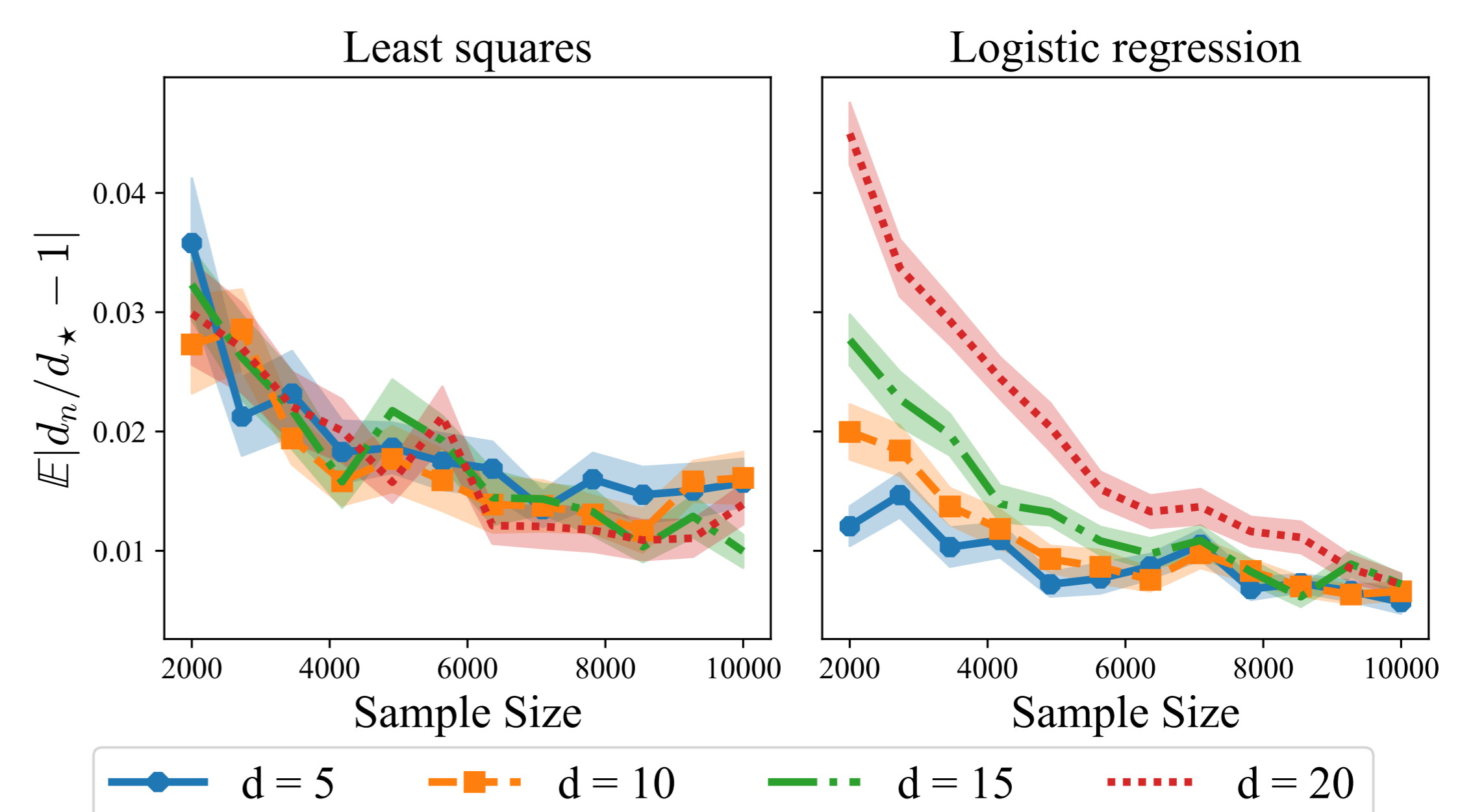
where $d_* := \text{Tr}(\Omega(\theta_*))$ is the **effective dimension**.

- When the model is **well-specified**, $d_* = d$ and $\|\Omega(\theta_*)\|_2 = 1$.
- When the model is **misspecified**,

	Eigendecay		Dimension Dependency		Ratio
	G_*	H_*	d_*	d	d_*/d
Poly-Poly	$i^{-\alpha}$	$i^{-\beta}$	$d^{(\beta-\alpha+1)\vee 0}$	d	$d^{(\beta-\alpha)\vee(-1)}$
Poly-Exp	$i^{-\alpha}$	$e^{-\nu i}$	$d^{1-\alpha} e^{\nu d}$	d	$d^{-\alpha} e^{\nu d}$
Exp-Poly	$e^{-\mu i}$	$i^{-\beta}$	1	d	d^{-1}
Exp-Exp	$e^{-\mu i}$	$e^{-\nu i}$	d if $\mu = \nu$	d	1 if $\mu = \nu$
			1 if $\mu > \nu$	d	d^{-1} if $\mu > \nu$
			$e^{(\nu-\mu)d}$ if $\mu < \nu$		$d^{-1} e^{(\nu-\mu)d}$ if $\mu < \nu$

The bound matches the **Cramér-Rao** lower bound.

- d_* can be approximated by $d_n := \text{Tr}(G_n(\hat{\theta}_n)^{1/2} H_n(\hat{\theta}_n)^{-1} G_n(\hat{\theta}_n)^{1/2})$.
- **How well does d_n approximate d_* ?**



Examples

Generalized linear models. Consider the statistical model

$$p_\theta(y | x) \sim \frac{\exp[\theta^\top t(x, y) + h(x, y)]}{\int \exp[\theta^\top t(x, \bar{y}) + h(x, \bar{y})] d\mu(\bar{y})} d\mu(y)$$

with $\|t(X, Y)\|_2 \leq_{a.s.} M$. It induces the loss function

$$\ell(\theta; z) := -\theta^\top t(x, y) - h(x, y) + \log \int \exp[\theta^\top t(x, \bar{y}) + h(x, \bar{y})] d\mu(\bar{y}),$$

which is $(2M, 2)$ -generalized self-concordant.

Score matching with exponential families. Consider an exponential family with density $\log p_\theta(z) = \theta^\top t(z) + h(z) - \Lambda(\theta)$. The score matching loss is

$$\ell(\theta; z) = \frac{1}{2} \theta^\top A(z) \theta - b(z)^\top \theta + c(z) + \text{const},$$

where $A(z) := \sum_{k=1}^p \frac{\partial t(z)}{\partial z_k} (\frac{\partial t(z)}{\partial z_k})^\top$ is positive semi-definite,

$$b(z) := \sum_{k=1}^p \left[\frac{\partial^2 t(z)}{\partial z_k^2} + \frac{\partial h(z)}{\partial z_k} \frac{\partial t(z)}{\partial z_k} \right], \quad \text{and} \quad c(z) := \sum_{k=1}^p \left[\frac{\partial^2 h(z)}{\partial z_k^2} + \left(\frac{\partial h(z)}{\partial z_k} \right)^2 \right].$$

It is generalized self-concordant for all $\nu \geq 2$ and $R \geq 0$.