



A Tutorial on the Non-Asymptotic Theory of System Identification

CDC'23 and <https://arxiv.org/abs/2309.03873>

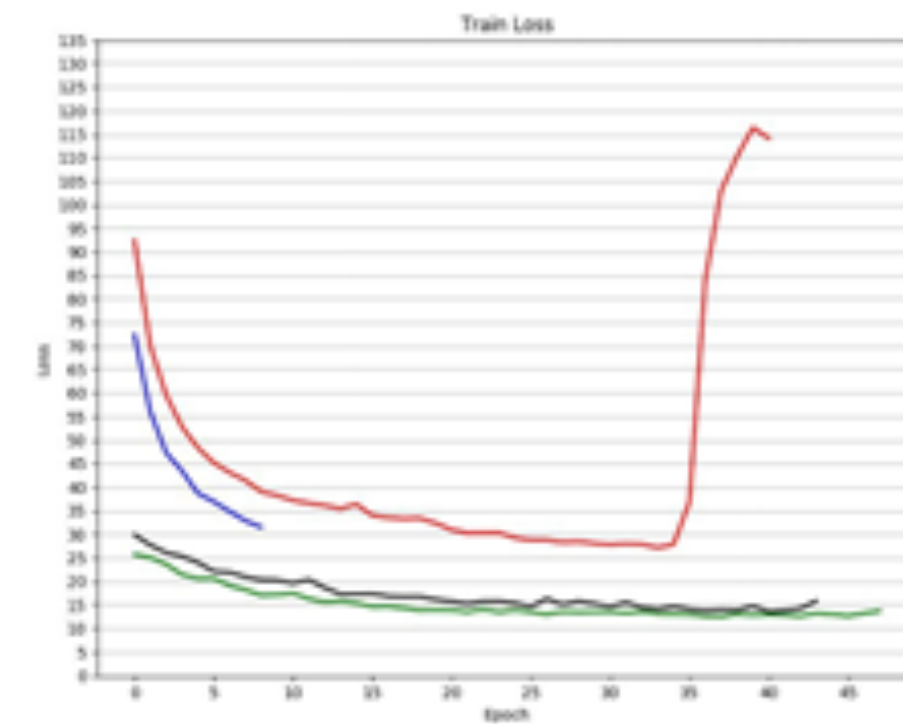
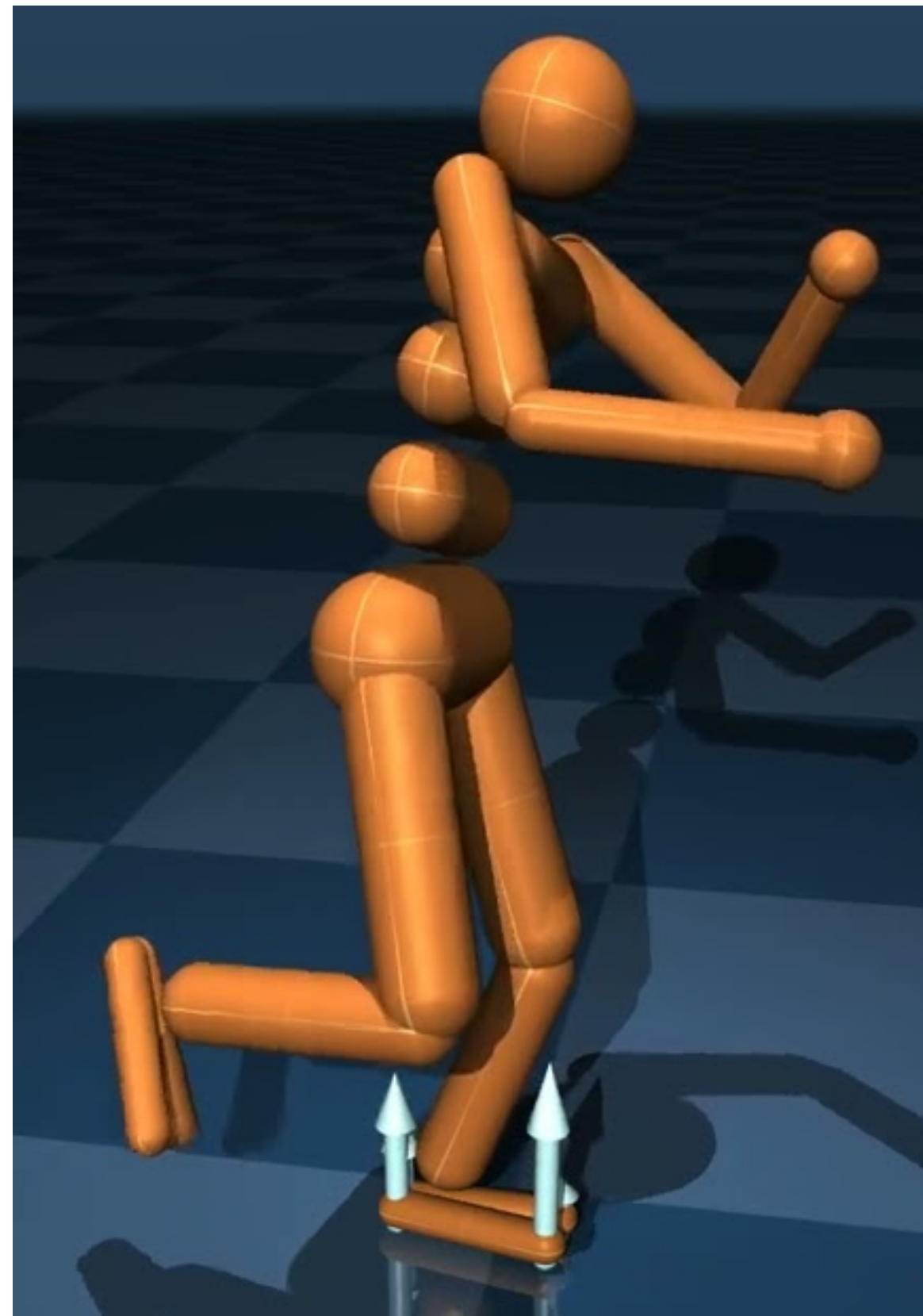


Ingvar Ziemann (Penn), Anastasios Tsiamis (ETH), Bruce Lee (Penn), Yassir Jedra (MIT), Nikolai Matni (Penn), George J. Pappas (Penn)

A snapshot of what lead us here



Ambition:



Not just in sim:



Reality:



A snapshot of what lead us here



Journals & Magazines > IEEE Control Systems Magazine > Volume: 43 Issue: 5

appears in part 2/2

Data-Driven Control: Part One of Two



Statistical Learning Theory for Control

A FINITE SAMPLE PERSPECTIVE

Anastasios Tsiamis*, Ingvar Ziemann*, Nikolai Matni, and George J. Pappas

Unknown System

$$X_{t+1} = A^* X_t + B^* U_t + W_t, \quad U_t = K X_t$$

Find K such that

$$\sum_{t=1}^T \mathbf{E} \{ X_t^\top Q X_t + U_t^\top R U_t \} = \min$$

today!

experiments

data

$$(X, U)_{1:T}$$

system identification

model estimates

$$(\hat{A}, \hat{B})$$

control

$$U_t = \hat{K} X_t$$

$$U_t \sim N(0, I)$$

<https://arxiv.org/abs/2209.05423>

Overview



Introduction and Roadmap (20 minutes)

Concentration Inequalities (30 minutes)

Hanson-Wright Inequality and Self-Normalized Martingales

Lower Tail of the Empirical Covariance (20 min)

Concentration \neq Persistence of Excitation

System Identification (30 min)

ARX-identification

An Alternate Approach: The Offset Basic Inequality (20 min)

Extension to nonlinear problems

Statistical Setup



Consider a time series model

$$Y_t = \theta^* X_t + V_t, \quad t = 1, \dots, T$$

Where:

benign noise

Y_t - Outputs in \mathbb{R}^{d_Y}

X_t - Covariates in \mathbb{R}^{d_X}

V_t - Noise in \mathbb{R}^{d_Y}

θ^* - Unknown Parameter in $\mathbb{R}^{d_Y \times d_X}$

Example ARX(p,q):

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + W_t$$

In other words...

$$X_t = \begin{bmatrix} Y_{t-1:t-p}^\top & U_{t-1:t-q}^\top \end{bmatrix}^\top$$

$$\theta^* = \begin{bmatrix} A_{1:p}^* & B_{1:q}^* \end{bmatrix}$$

$$V_t = W_t.$$

Least Squares Estimation (LSE)



Consider a time-series model:

$$Y_t = \theta^* X_t + V_t, \quad t = 1, \dots, T$$

Least Squares Estimator:

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^{d_Y \times d_X}} \left\{ \frac{1}{T} \sum_{t=1}^T \|Y_t - \theta X_t\|_2^2 \right\}$$

\Rightarrow

$$\hat{\theta} \triangleq \left(\sum_{t=1}^T Y_t X_t^\top \right) \left(\sum_{t=1}^T X_t X_t^\top \right)^{-1}$$

Interested in:

$$\hat{\theta} - \theta^* = \left(\sum_{t=1}^T V_t X_t^\top \right) \left(\sum_{t=1}^T X_t X_t^\top \right)^{-1}$$

Today: Modern perspective on LSE

Draw on tools from:

Machine Learning Theory

High-Dimensional Statistics

High-Dimensional Probability

Problem

Establish finite sample guarantees:



$$\|\hat{\theta} - \theta^*\| \leq \epsilon \quad \text{wpa.} \quad 1 - \delta$$

Fix:

accuracy $\epsilon > 0$

failure probability $\delta \in (0,1)$

a norm $\|\cdot\|$

and a ‘reasonable’ estimator $\hat{\theta}$

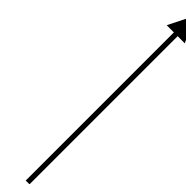
Typically we can prove:

$$\epsilon \propto (\text{noise scale}) \times \sqrt{\frac{\text{dimension} + \log(1/\delta)}{\text{sample size}}}$$

As long as:

$$\text{sample size} \gtrsim \text{dimension} + \log(1/\delta)$$

Persistence of Excitation



The Path Ahead

Covariance Concentration



At LLN Scale ~ 1

$$\hat{\theta} - \theta^* = \left(\sum_{t=1}^T V_t X_t^\top \right) \left(\sum_{t=1}^T X_t X_t^\top \right)^{-1}$$
$$\hat{\theta} - \theta^* = \frac{1}{T} \left[\left(\sum_{t=1}^T V_t X_t^\top \right) \left(\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right] \left(\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \right)^{-1/2}$$

← Random Matrix!

Todo:

Random Walk at CLT Scale $\sim \sqrt{T}$

CLT-analogue: Self-Normalized Martingale Ineq.

LLN-analogue: Covariance (anti-)Concentration



Finite sample aspects and comparison with asymptotic

Finite Sample Guarantees

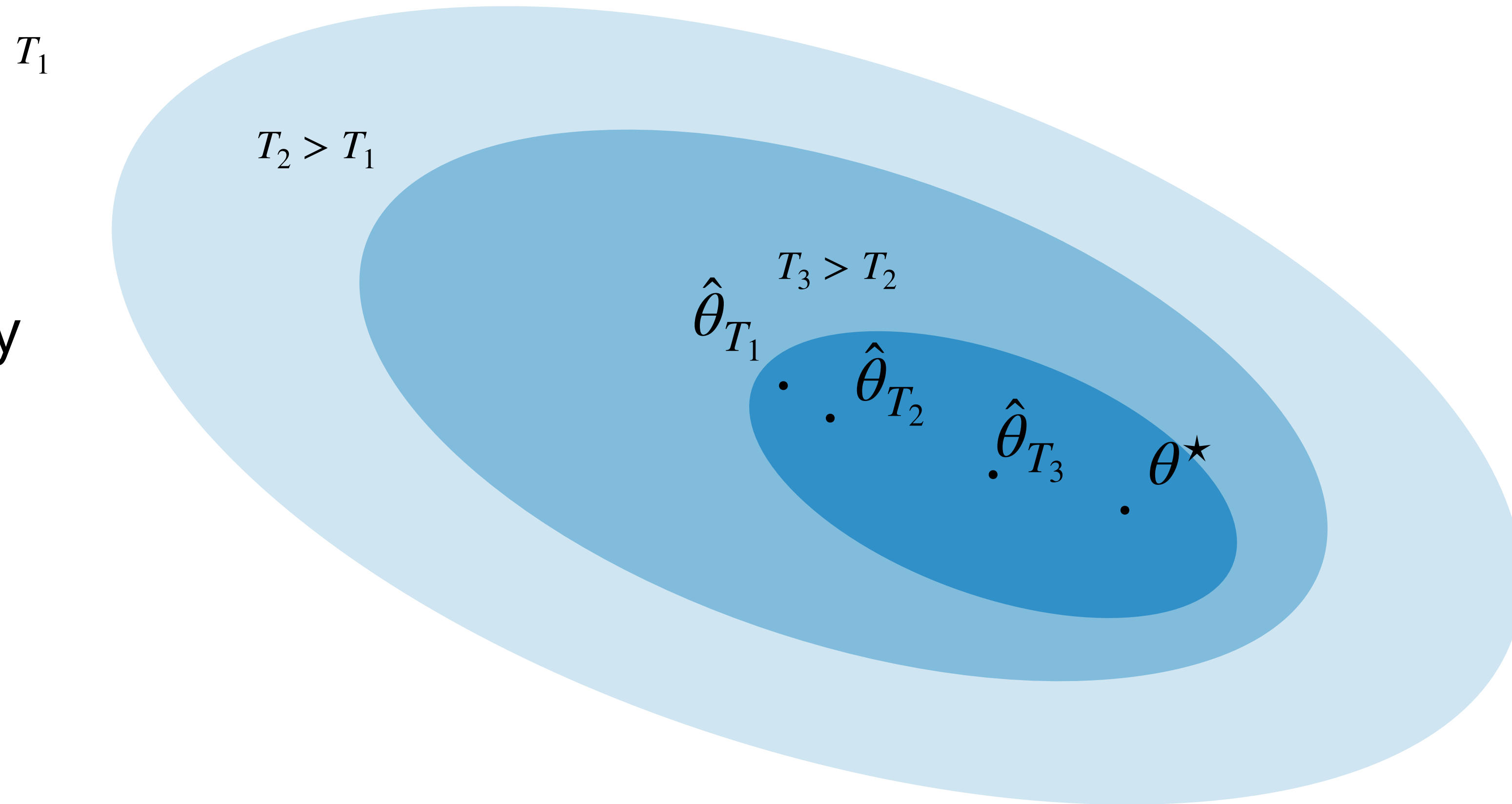


What is the error $\hat{\theta}_T - \theta^*$?

Finite Sample Guarantees

$$\theta^* \in \mathcal{B}(T)$$

Q1: Error decay rate?

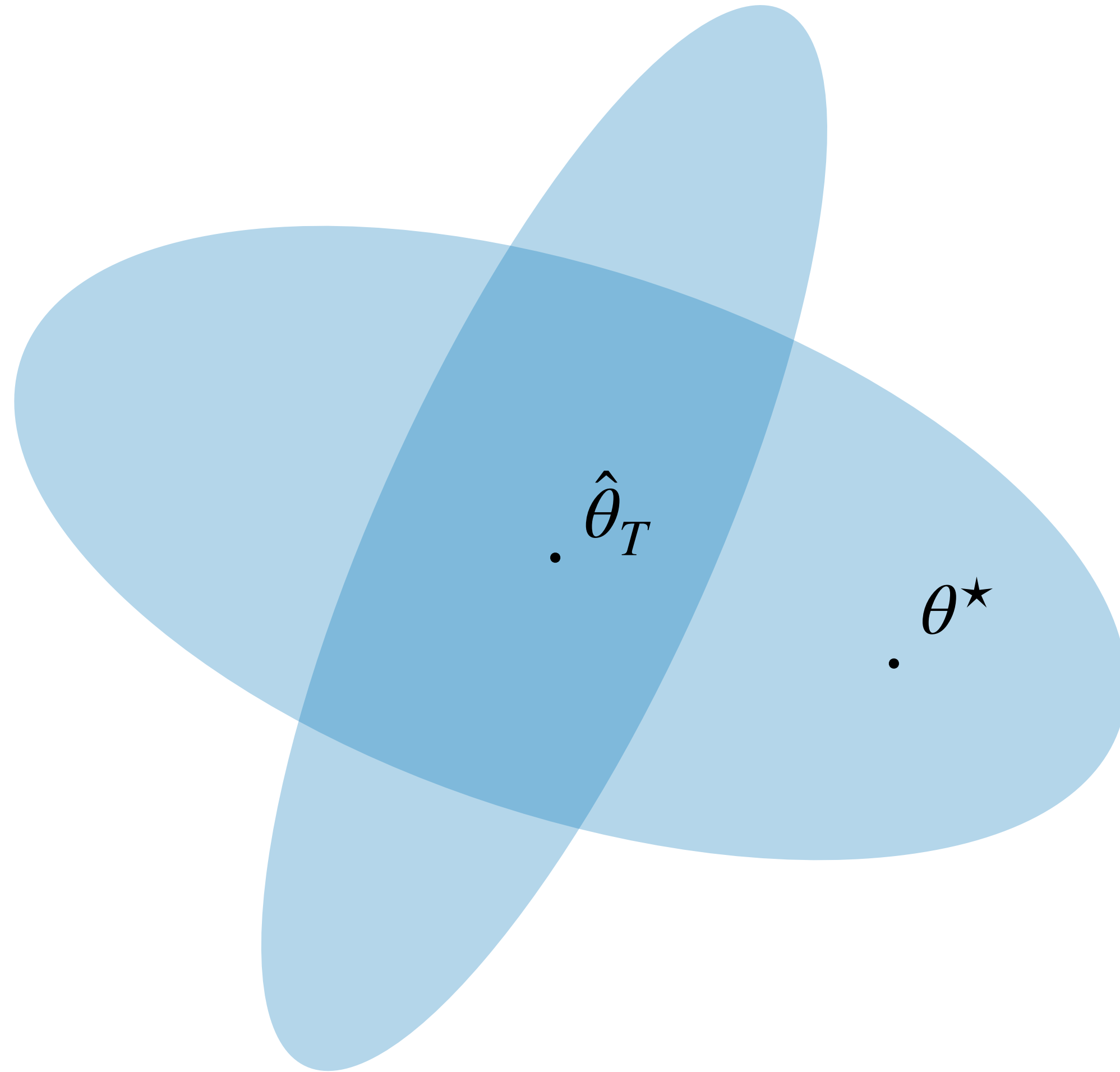


Finite Sample Guarantees

$$\theta^* \in \mathcal{B}(T)$$

Q1: Error decay rate?

Q2: Shape of uncertainty?



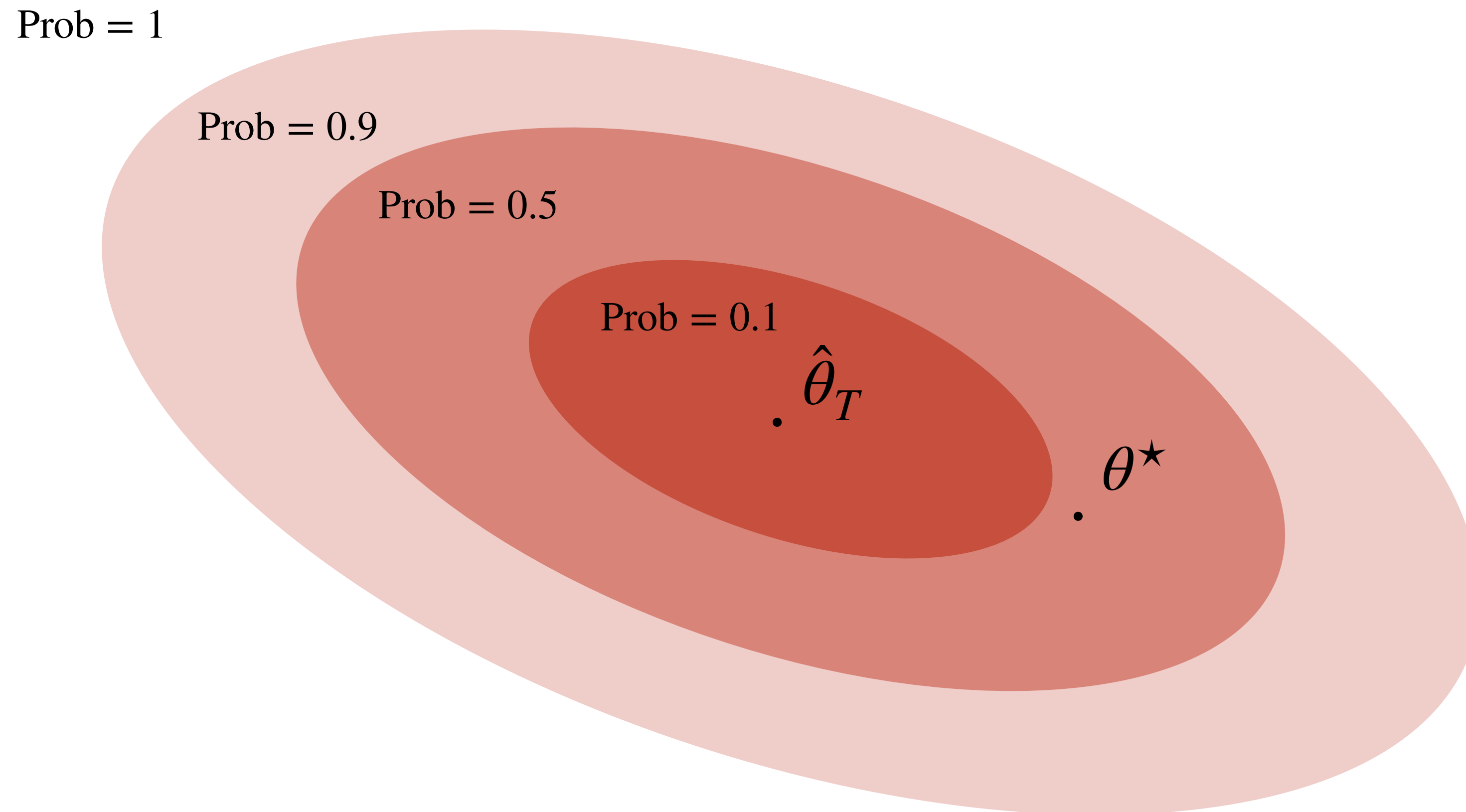
Finite Sample Guarantees

$$\text{Prob}(\theta^* \in \mathcal{B}(\hat{\theta}_T, \delta)) \geq 1 - \delta$$

Q1: Error decay rate?

Q2: Shape of uncertainty?

Q3: Confidence?



Finite Sample Guarantees

$$\mathbb{P}(\theta^* \in \mathcal{B}(c, \delta)) \geq 1 - \delta$$

System specific
& Universal constants

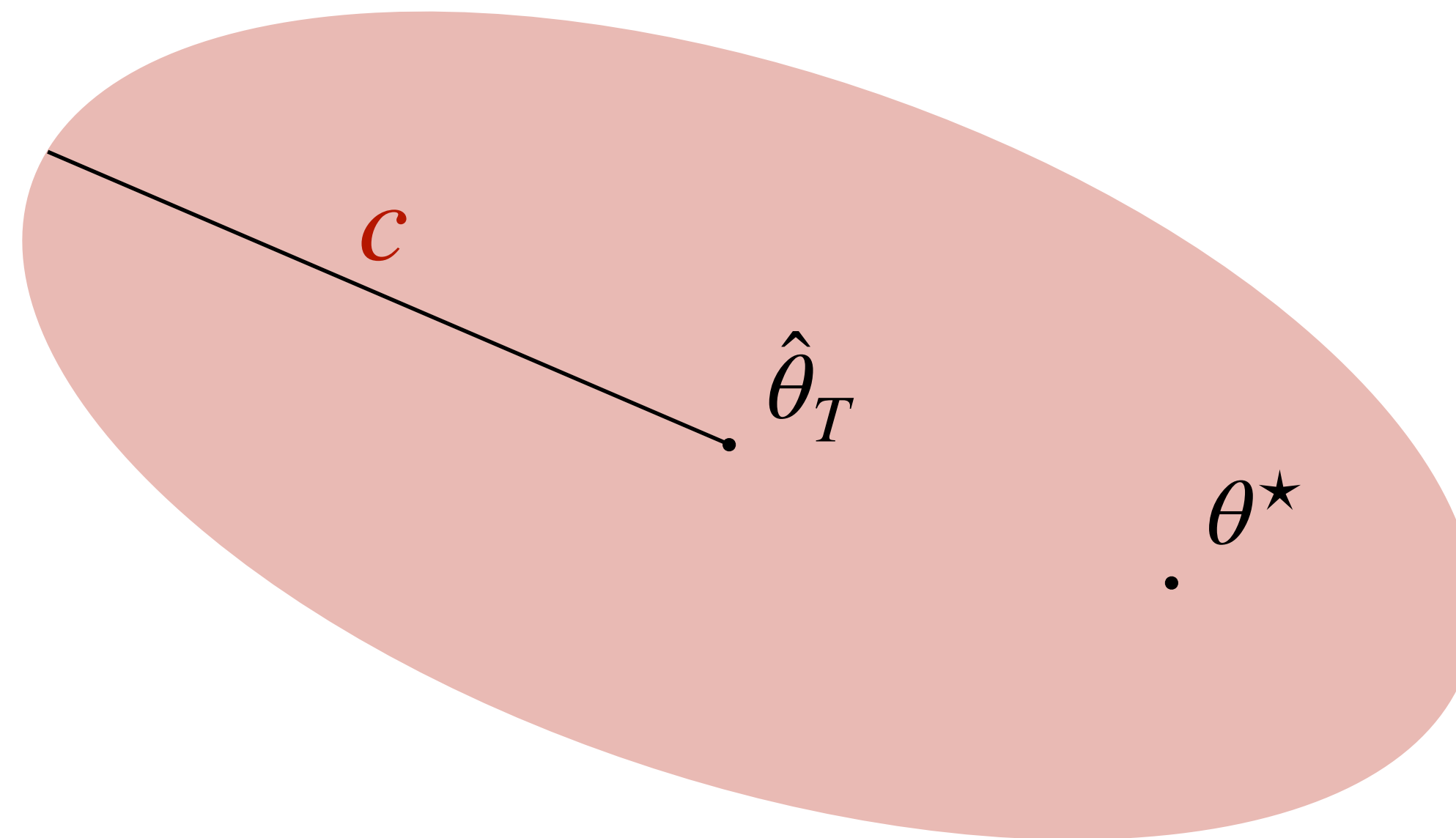
Relative scaling

Q1: Error decay rate?

Q2: Shape of uncertainty?

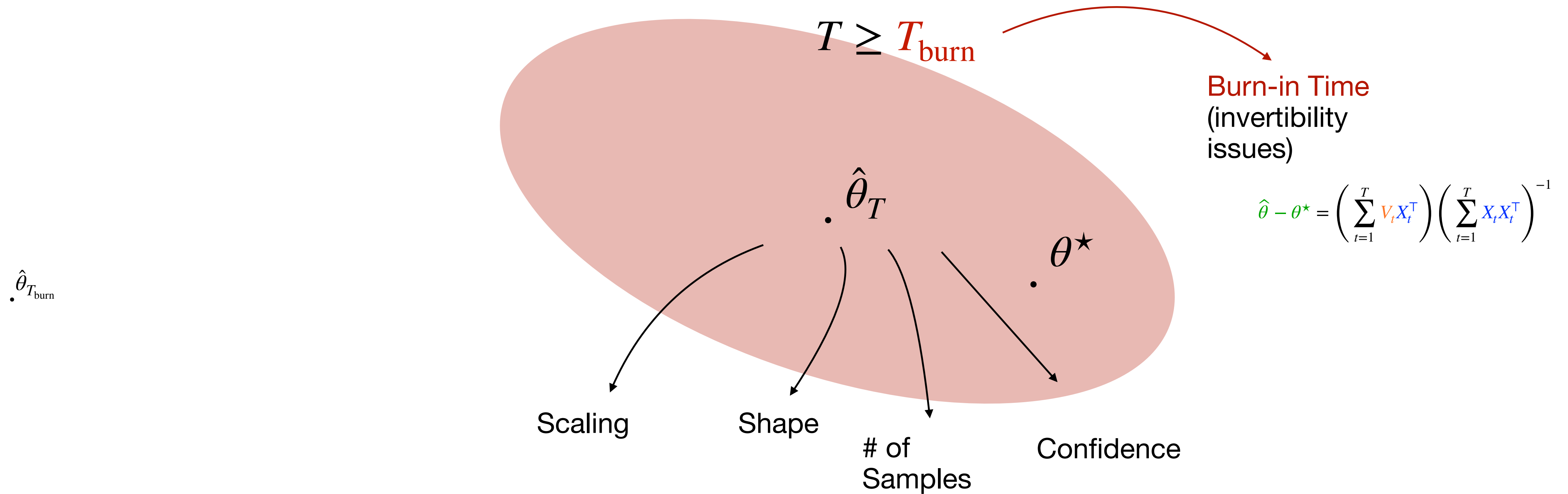
Q3: Confidence?

Q4: Absolute Scaling?



Finite Sample Guarantees

$$\text{Prob}(\theta^* \in c\mathcal{B}(T, \delta)) \geq 1 - \delta$$



Asymptotics

$$\sqrt{T}(\hat{\theta}_T - \theta^*) \xrightarrow{d} \mathcal{N}(0, \Sigma_\theta)$$

Ljung 1999

Deviation for
 $T < \infty$?

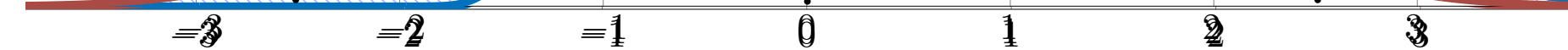
Using finite samples
 $\exp(-T)$

Burn-in
 $1/\delta^2$

Biggest Difference
(Berry Esseen)
 $1/\sqrt{T}$

Burn-in
 $\log 1/\delta$

However negative
events: **tail
events**
Seem to decay
much faster



Why finite sample?

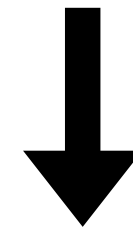
Complementary
tools

	Finite Sample	Asymptotics
Rate	$1/\sqrt{T}$	$1/\sqrt{T}, T \rightarrow \infty$
Shape	$\frac{1}{T} \sum_{t=1}^T X_t X_t^\top$	$\lim \frac{1}{T} \sum_{t=1}^T X_t X_t^\top$
Confidence	$T \geq \log 1/\delta$	$T \geq 1/\delta^2$
Scale	Conservative universal const.	Optimal
Transient	Burn-in times	×

With Berry Esseen

Renewed Attention

$$\|\hat{\theta}_T - \theta^*\| = O(1/\sqrt{T})$$



$$\|\hat{\theta}_T - \theta^*\| = \frac{C_{\text{sys}}}{\sqrt{T}}, \quad T \geq T_{\text{burn}}$$

System theoretic
properties

What makes learning difficult or hard: $C_{\text{sys}}, T_{\text{burn}}$



Thank you!