

Data-Driven Algorithm Design and Verification for Parametric Convex Optimization

Bartolomeo Stellato

Department of Operations Research and Financial Engineering

Department of Electrical and Computer Engineering

Department of Computer Science







**Most applications require fast and effective
decisions in real-time**

Real-time optimization can help us

$$\begin{array}{l} \text{minimize} \quad f(z, x) \\ \text{subject to} \quad z \in C(x) \end{array}$$

↑
parameters

↓
decisions

objective f : energy consumption, costs

constraints C : dynamics, physical limits

re-planning in real-time
is the key to effective
decision-making

How do we solve
such problems?

...and they can solve many constrained convex problems!

Linear Programs



PDLP

Applegate, Díaz, Hinder, Lu, Lubin,
O'Donoghue, Schudy (2021)

Quadratic Programs



OSQP

Stellato, Banjac, Goulart,
Bemporad, Boyd (2020)

Conic Programs



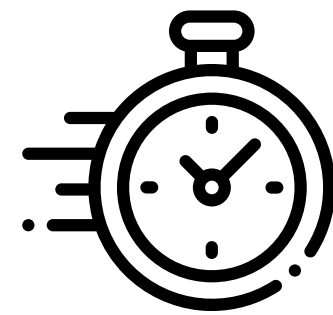
SCS
SPLITTING CONIC SOLVER

O'Donoghue, Chu, Parikh, Boyd (2016)

But they can converge slowly

major issue in safety-critical applications with

real-time
requirements



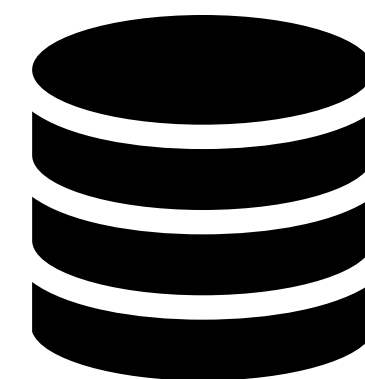
limited
computing power



💡 main idea

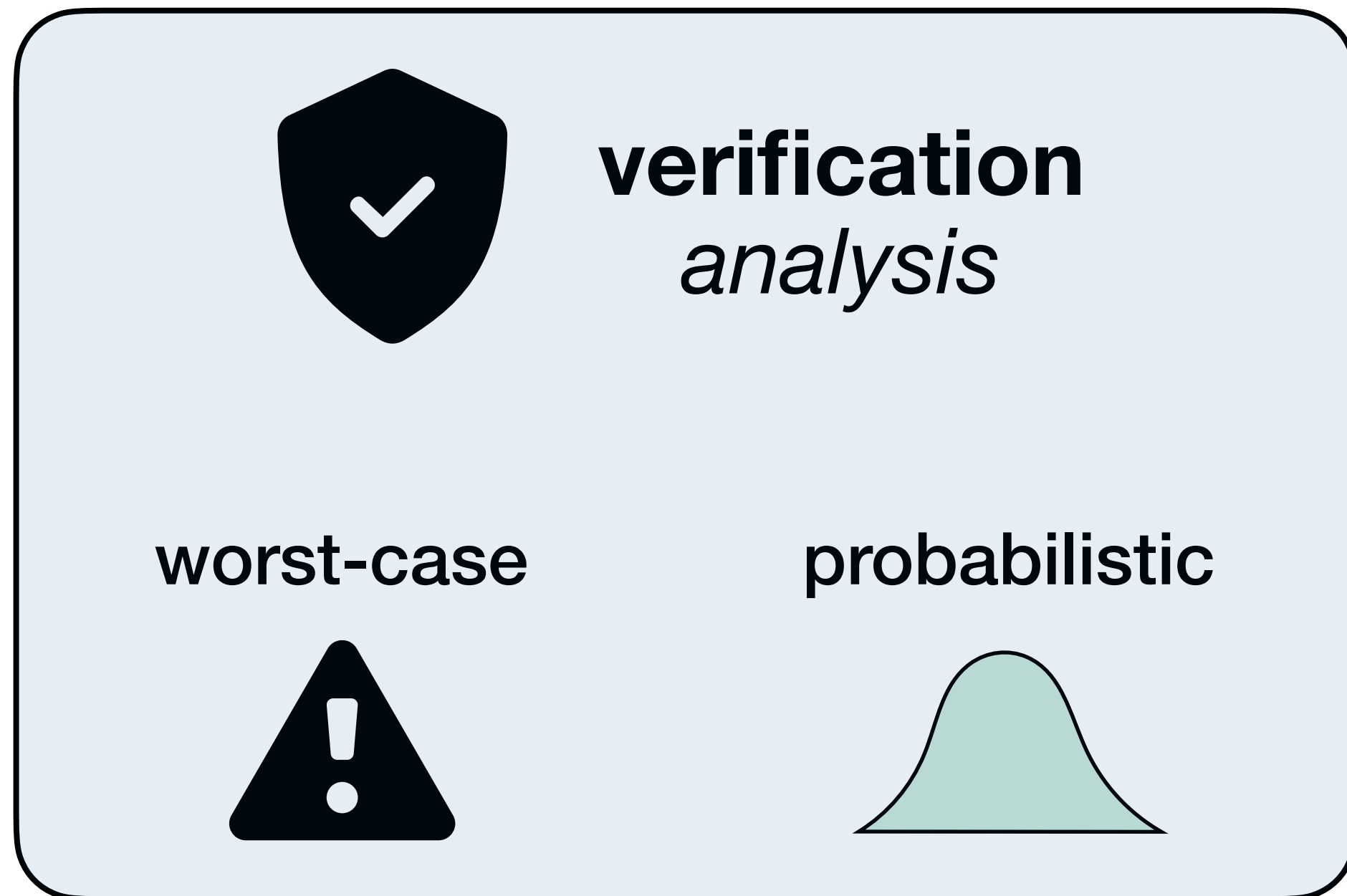
in most applications we repeatedly
solve the **same problem** with
varying parameters

minimize $f(z, \mathbf{x})$
subject to $z \in C(\mathbf{x})$



large amount of data
(e.g., instances, solutions)

First-order methods in parametric convex optimization



A light blue rounded rectangle containing icons and text for verification analysis. At the top left is a shield with a checkmark. To its right is the text "verification analysis". Below this, "worst-case" is written above a warning triangle icon, and "probabilistic" is written above a bell curve icon.

verification
analysis

worst-case probabilistic



A white rounded rectangle containing icons and text for design learning. At the top left is a gear icon. To its right is the text "design learning". Below this is the text "with probabilistic guarantees" above a bell curve icon that contains a shield with a checkmark.

design
learning

with probabilistic
guarantees

Performance verification

Convergence of first-order methods

iterations

$$z^{k+1} = T(z^k, x) \quad \text{for } k = 0, 1, \dots$$

T operator
(e.g., contractive, averaged)

goal: find fixed-points

$$z^* = T(z^*, x)$$

example
gradient descent

problem \longrightarrow optimality conditions
minimize $f(z, x) \longrightarrow \nabla f(z^*, x) = 0$

iterations
$$z^{k+1} = z^k - \theta \nabla f(z^k, x)$$

fixed-points
$$z^* = z^* - \theta \nabla f(z^*, x)$$

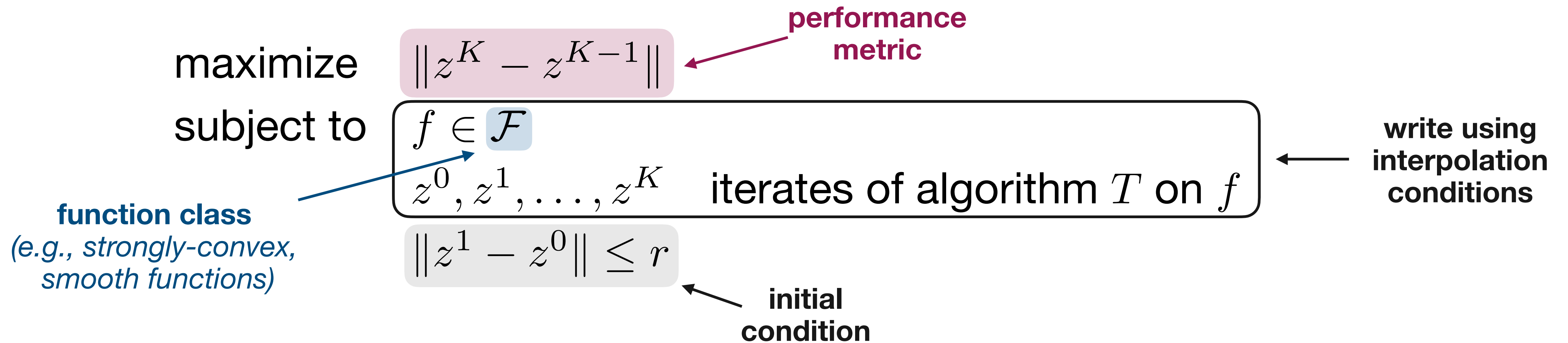
same as

performance metric

$$r^k(x) = \|T(z^{k-1}) - z^{k-1}\| = \|z^k - z^{k-1}\|$$

fixed-point residual
(converges to 0)

Classical convergence bounds via Performance Estimation



convex SDP
Gram matrix reformulation

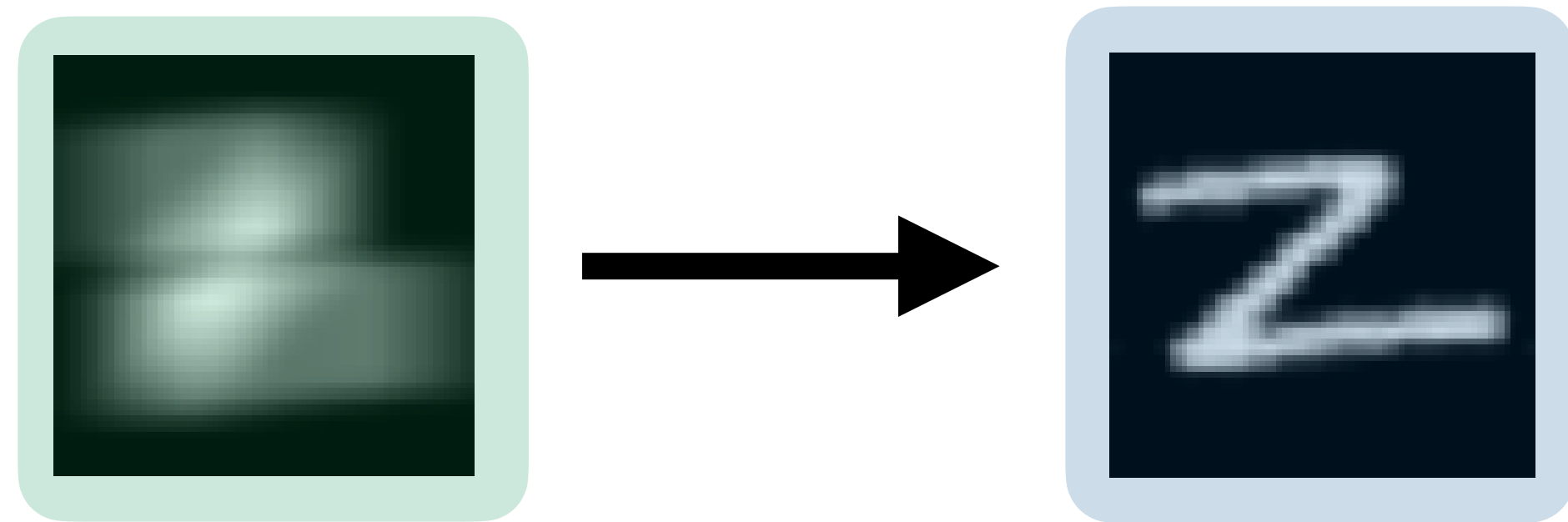
$$G = \begin{bmatrix} \|z^1 - z^0\|_2^2 & (z^1 - z^0)^T g^1 & (z^1 - z^0)^T g^0 & \dots \\ (z^1 - z^0)^T g^1 & \|g^1\|_2^2 & (g^1)^T g^0 & \dots \\ (z^1 - z^0)^T g^0 & (g^1)^T g^0 & \|g^0\|_2^2 & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix}$$

gradients

independent from
iterate dimensions

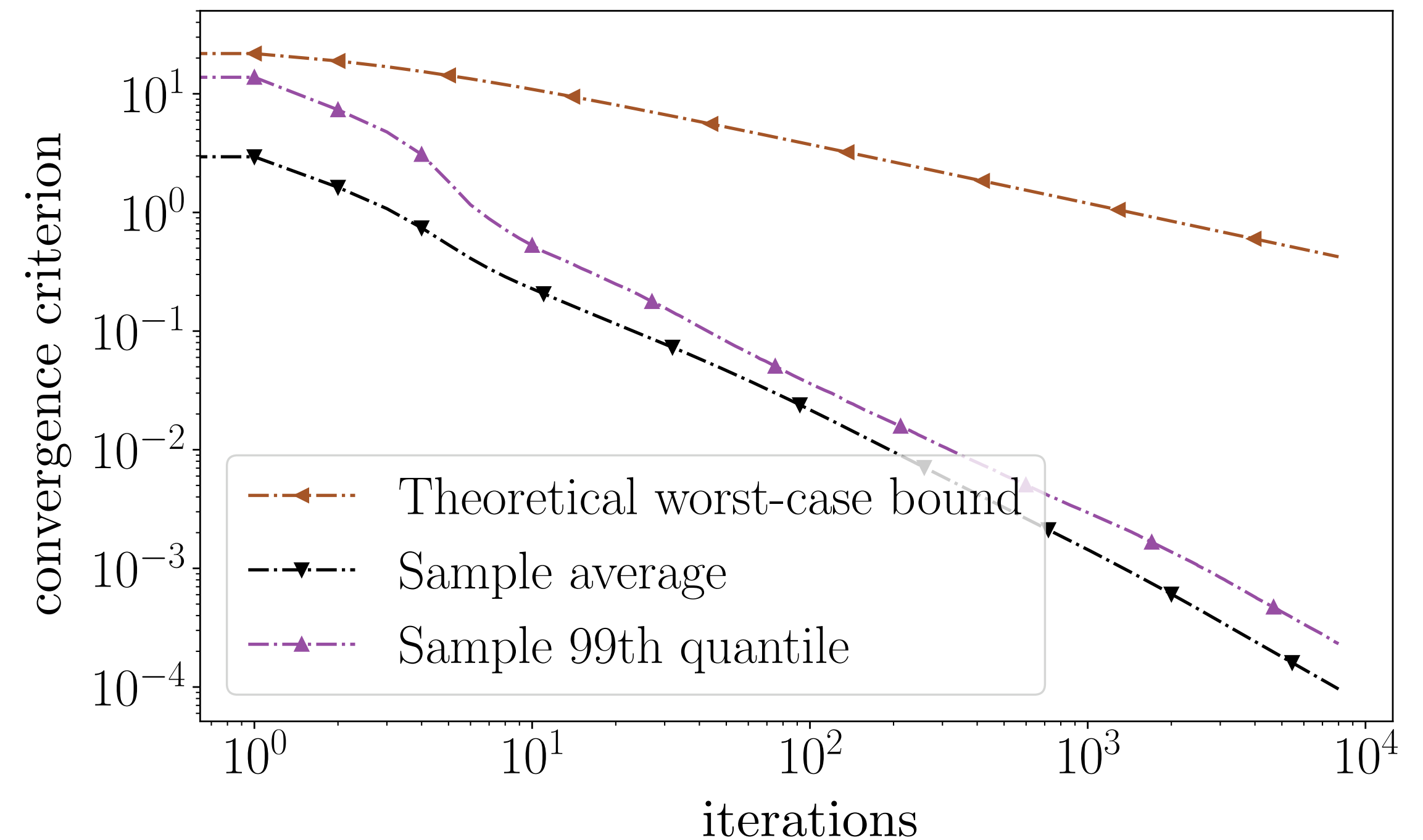
Classical worst-case convergence bounds can be very loose

image deblurring problem
emnist dataset



minimize $\|Az - x\|_2^2 + \lambda \|z\|_1$
subject to $0 \leq z \leq 1$

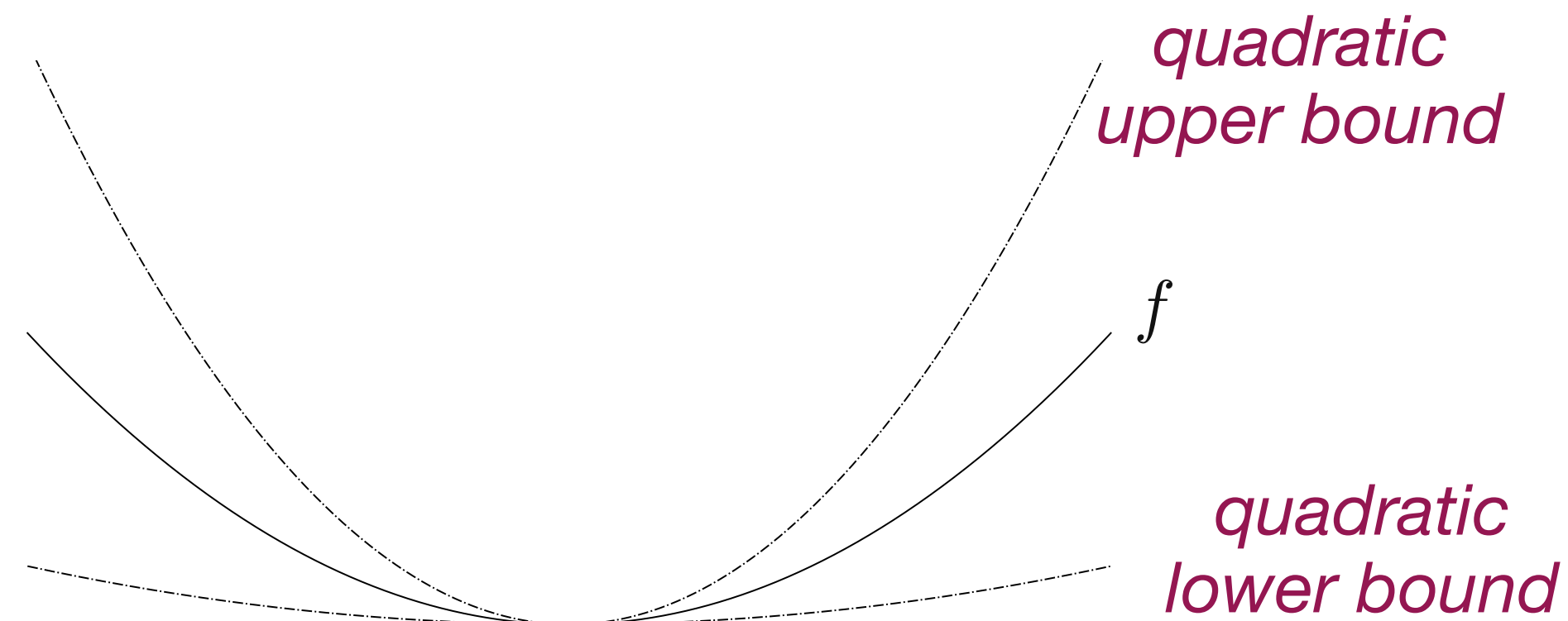
deblurred image (pointing to z)
blurred image (pointing to x)



why are worst-case bounds pessimistic?

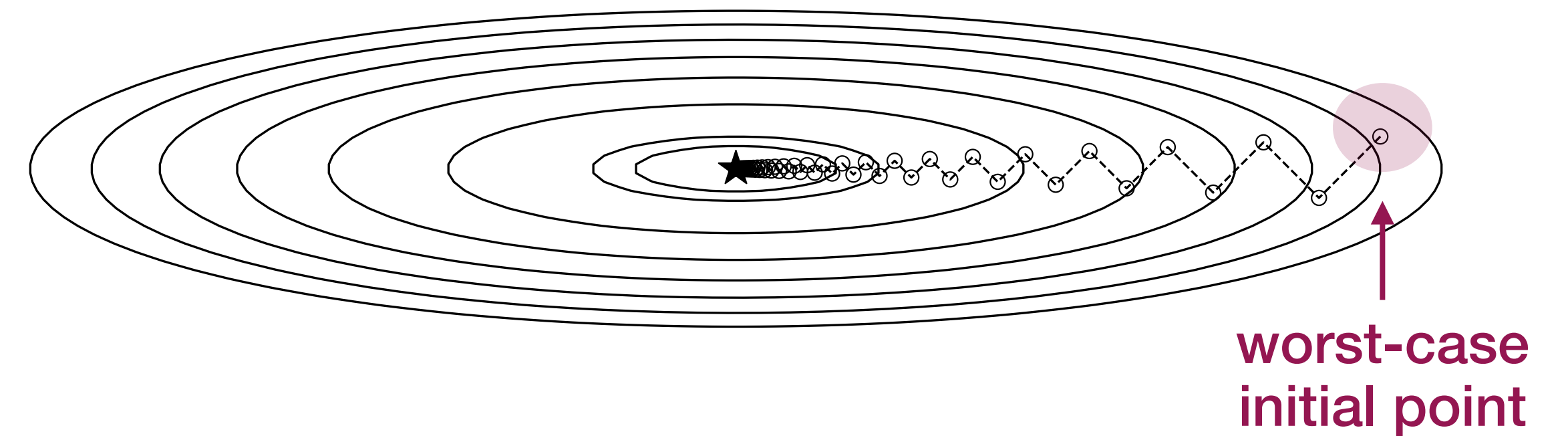
Issues with classical convergence analysis

general function classes
(f is strongly convex and smooth...)



we may never encounter that function

pessimistic bounds



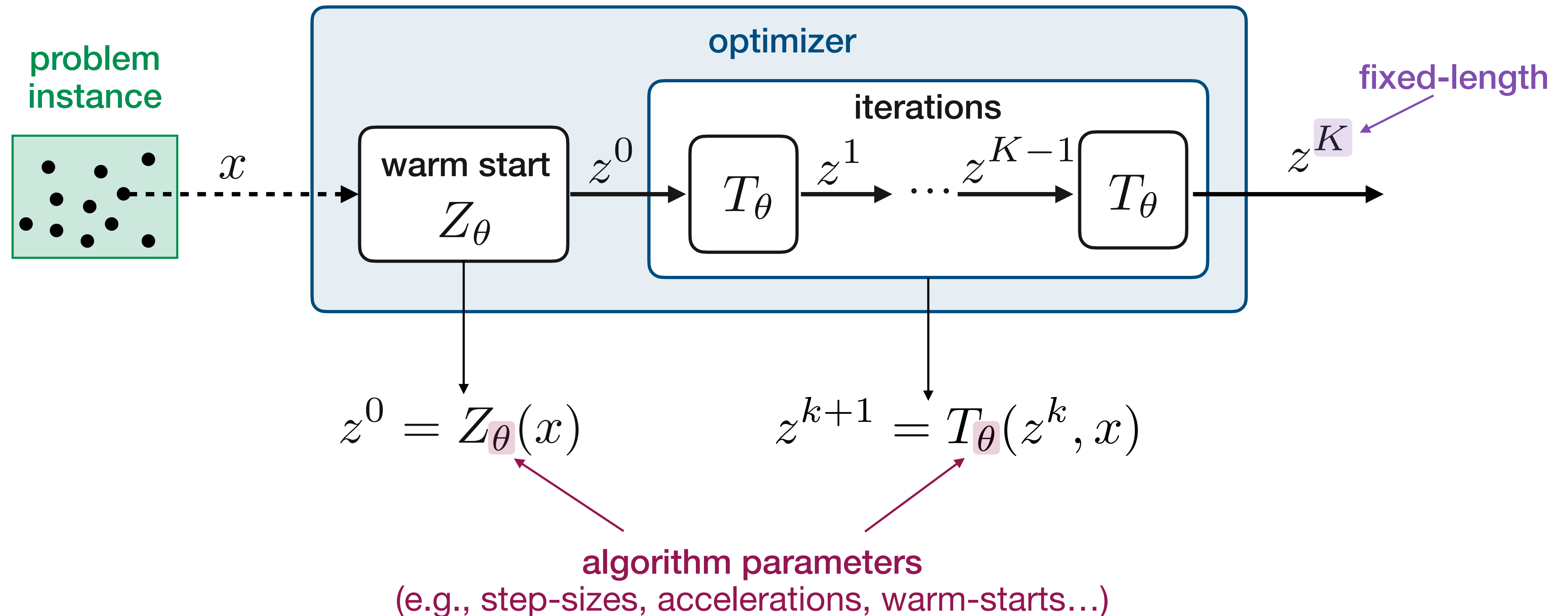
we may never start from that point

practical settings

$$\begin{aligned} &\text{minimize} && f(z, \mathbf{x}) \\ &\text{subject to} && z \in C(\mathbf{x}) \end{aligned}$$

same problem with varying parameters
 $\longrightarrow x \sim \mathbf{P}$
(unknown distribution)

Algorithms as fixed-length computational graphs



example
projected gradient descent

$$z^{k+1} = \Pi_{C(x)}(z^k - \theta \nabla_z f(z^k, x))$$

Verifying the algorithm performance after K iterations

goal

estimate norm of fixed-point residual

$$r^K(x) = \|z^K - z^{K-1}\|$$

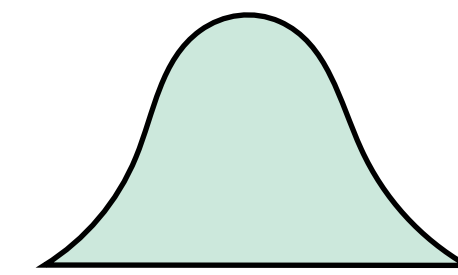


worst-case

$$\max_{x \in \mathcal{X}} r^K(x) \leq \epsilon$$

problem instances

convergence tolerance



probabilistic

$$\mathbf{P}(r^K(x) > \epsilon) \leq \eta$$

problem instances

convergence tolerance

probability bound



Worst-case algorithm verification

Parametric quadratic optimization

$$\max_{x \in \mathcal{X}} r^K(x) = \text{maximize } \|z^K - z^{K-1}\|$$

performance metric

$$\text{subject to } z^{k+1} = T_\theta(z^k, x), \quad k = 0, \dots, K-1$$

$$z^0 = Z_\theta(x), \quad x \in \mathcal{X}$$

problem instances

directly encode proximal algorithms
without interpolation inequalities

	step	verification constraint
affine (e.g., gradient, restarts, linear system solves)	$Dz^{k+1} = Az^k + Bx$	$Dz^{k+1} = Az^k + Bx$
elementwise maximum (e.g., separable projections, soft-thresholding,...)	$z^{k+1} = \max\{z^k, 0\}$	$z^{k+1} \geq 0, \quad z^{k+1} \geq z^k$ $(z^{k+1})^T (z^{k+1} - z^k) = 0$

similar to ReLU

similar to neural
network verification

Liu et al. (2021), Albarghouthi (2021)

Relaxing verification problem to an SDP

The verification problem is NP-hard
(by reduction from 0-1 integer programming)



convex semidefinite
 program relaxation

step	verification constraint	relaxed constraint
elementwise maximum <i>(e.g., box projections, soft-thresholding,...)</i>	$z^{k+1} = \max\{z^k, 0\}$	$z^{k+1} \geq 0, \quad z^{k+1} \geq z^k$ $\text{tr} \left(\begin{bmatrix} I & -I/2 \\ -I/2 & 0 \end{bmatrix} M \right) = 0$
	$z^{k+1} \geq 0, \quad z^{k+1} \geq z^k$ $(z^{k+1})^T (z^{k+1} - z^k) = 0$	$M \succeq \begin{bmatrix} z^{k+1} \\ z^k \end{bmatrix} \begin{bmatrix} z^{k+1} \\ z^k \end{bmatrix}^T$

depends on
 iterate dimensions

Unconstrained QP

Exact SDP reformulation

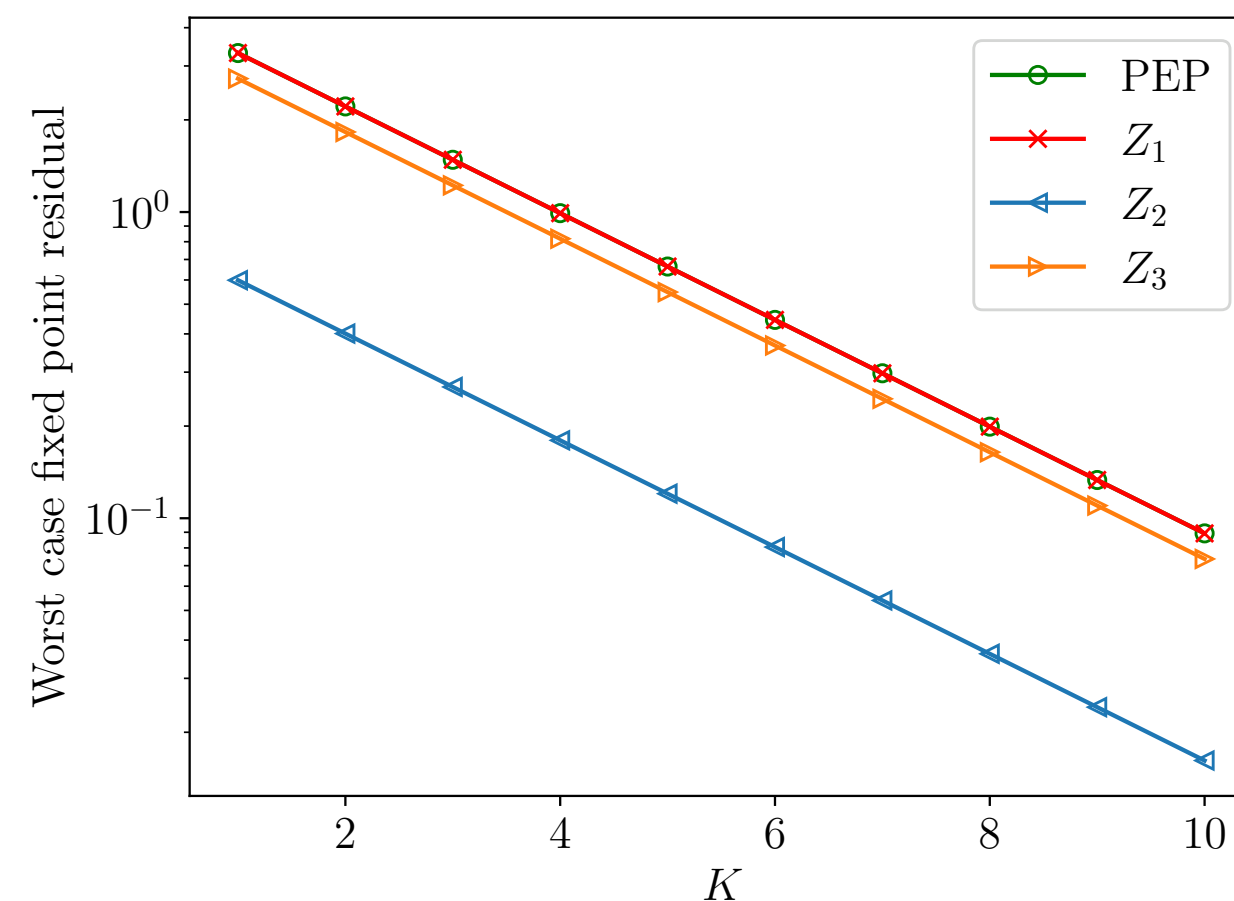
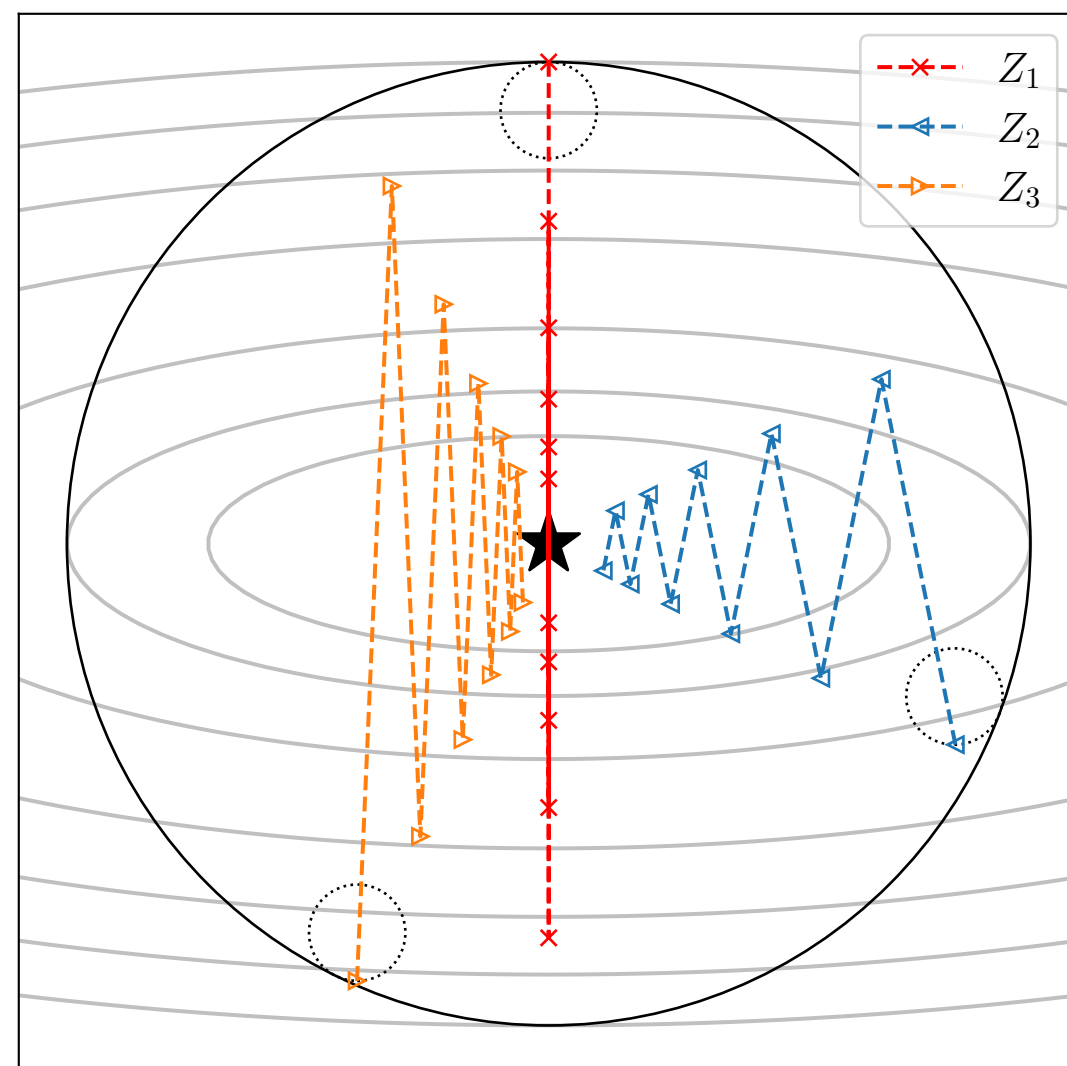
minimize $(1/2)z^T Pz + x^T z$ ← parameters

verification problem

maximize $\|z^K - z^{K-1}\|$
 subject to $z^{k+1} = z^k - \theta(Pz^k + x), \quad k = 0, \dots, K-1$ ← gradient descent
 $z^0 = Z_\theta(x), \quad x \in \mathcal{X}$

warm-starts

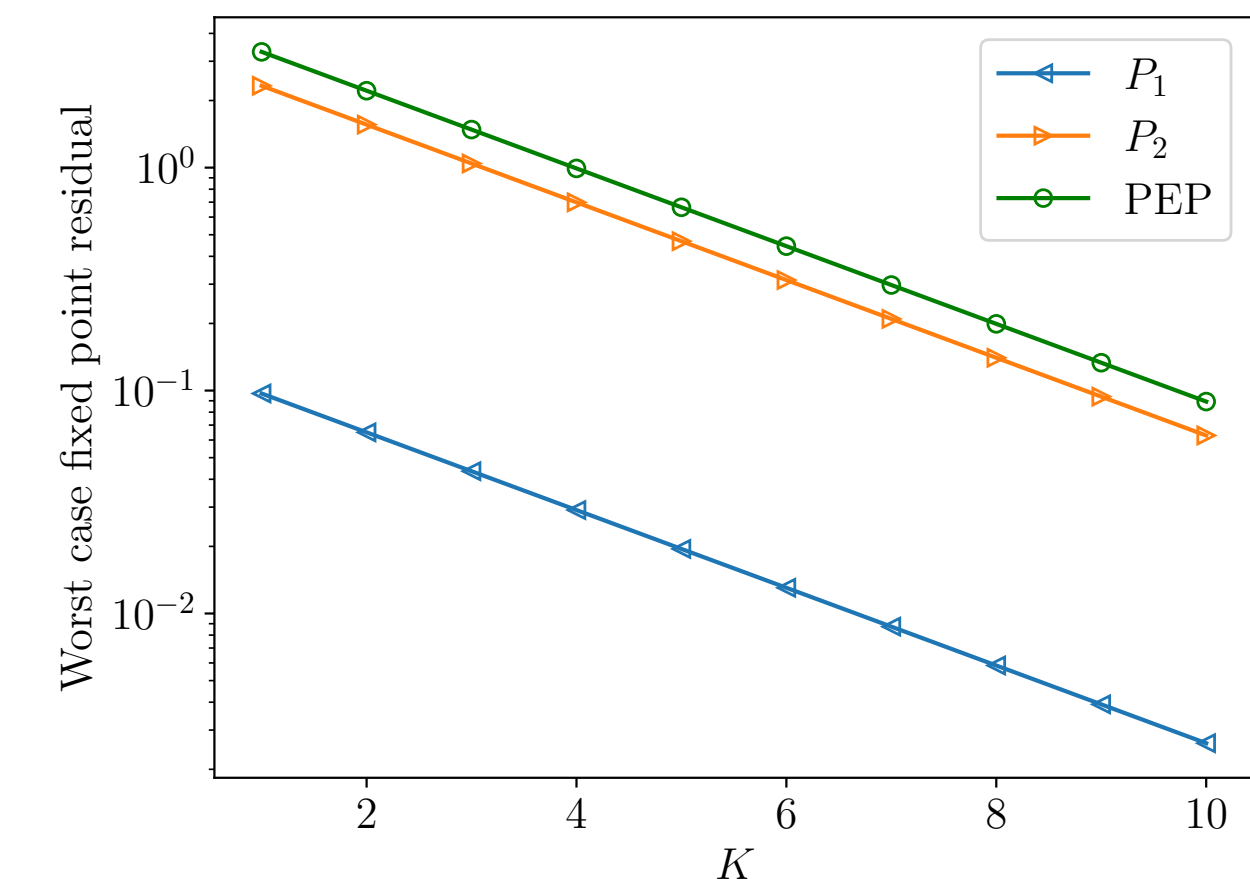
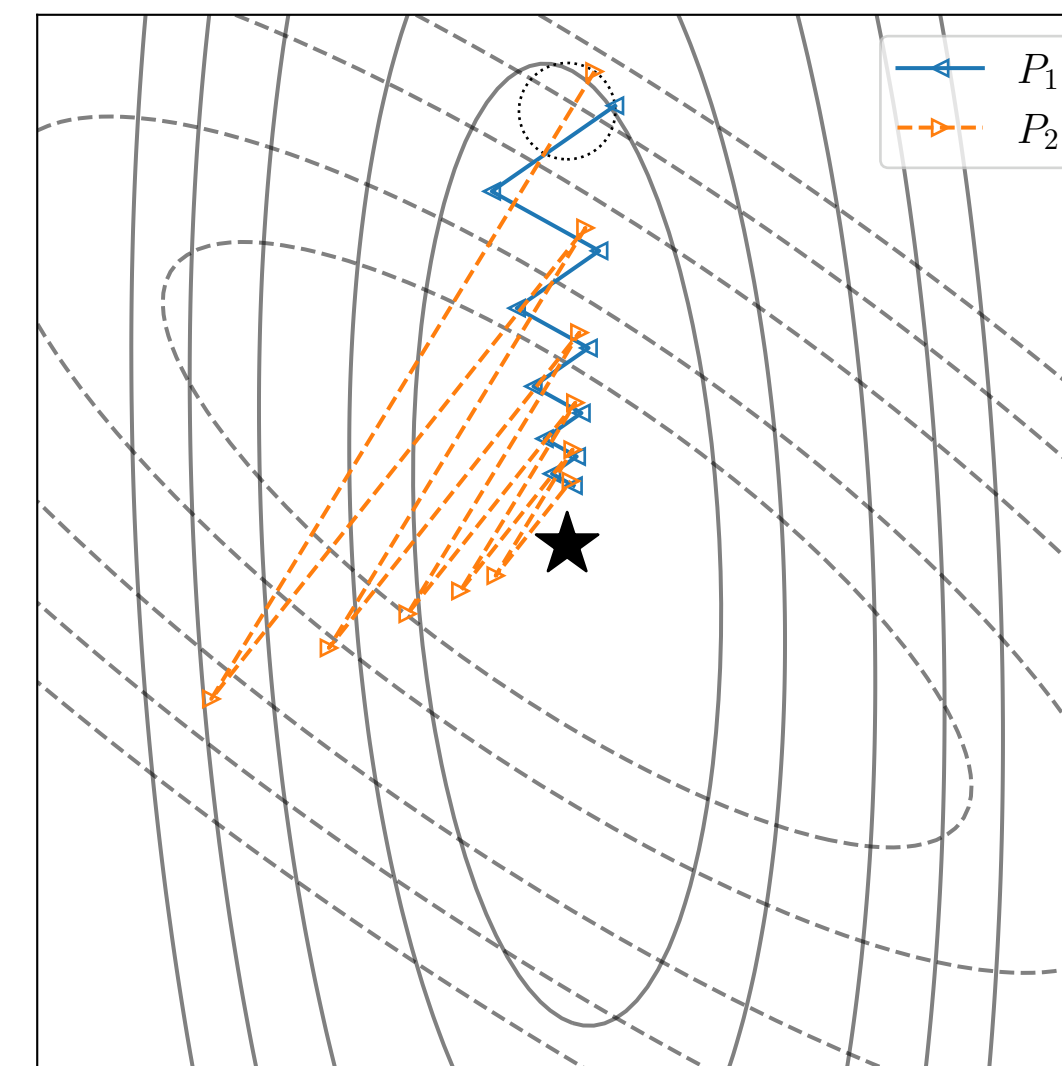
case I
 $Z_\theta(x) = Z_1, Z_2, \text{ or } Z_3$
 $x \in \mathcal{X} = \{0\}$



PEP cannot distinguish warm-starts

rotated functions

case II
 $Z_\theta(x) = \{z \mid \|z - 0.9 \cdot \mathbf{1}\| \leq 0.1\}$
 $x \in \mathcal{X} = \{0\}$
 P_1, P_2 rotations of P



PEP cannot distinguish quadratic functions

Nonnegative least-squares verification

nonnegative least squares

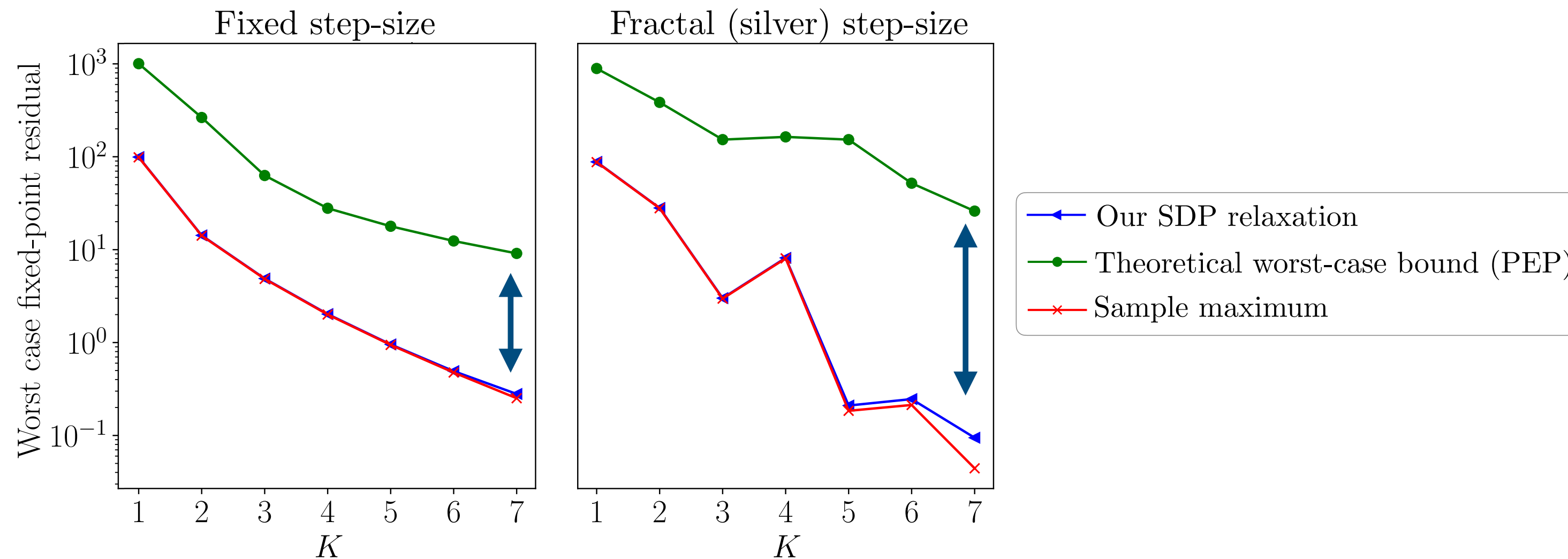
$$\begin{aligned} &\text{minimize} && (1/2) \|Az - x\|_2^2 \\ &\text{subject to} && z \geq 0 \end{aligned}$$

↑
parameters

verification problem

$$\begin{aligned} &\text{maximize} && \|z^K - z^{K-1}\| \\ &\text{subject to} && z^{k+1} = \max\{(I - \theta A^T A)z^k + \theta(A^T x), 0\}, \quad k = 0, \dots, K-1 \\ &&& z^0 = \{0\}, \quad x \in \mathcal{X} = \{x \mid \|x - 30 \cdot \mathbf{1}\| \leq 0.5\} \end{aligned}$$

projected
gradient
descent



10x-1000x reduction
(exploiting parametric structure)

computationally more expensive than PEP
(up to 1000 seconds for these instances)

Verification of First-Order Methods for Parametric Quadratic Optimization

V. Ranjan and B. Stellato

arXiv e-prints:2403.03331 (2024)

github.com/stellatogrp/algorithm_verification

Verifying the algorithm performance after K iterations

goal

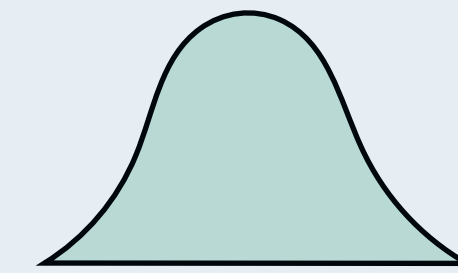
estimate norm of fixed-point residual

$$r^K(x) = \|z^K - z^{K-1}\|$$



worst-case

problem instances \rightarrow $\max_{x \in \mathcal{X}} r^K(x) \leq \epsilon$ \leftarrow convergence tolerance



probabilistic

problem instances \rightarrow $\mathbf{P}(r^K(x) > \epsilon) \leq \eta$ \leftarrow convergence tolerance \leftarrow probability bound

Probabilistic analysis

goal

*estimate probability of
computing bad-quality solutions*

$$\mathbf{P}(r^K(x) > \epsilon)$$

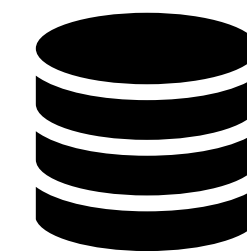
↑
any metric
(e.g., fixed-point residual)

issue

we don't know P !



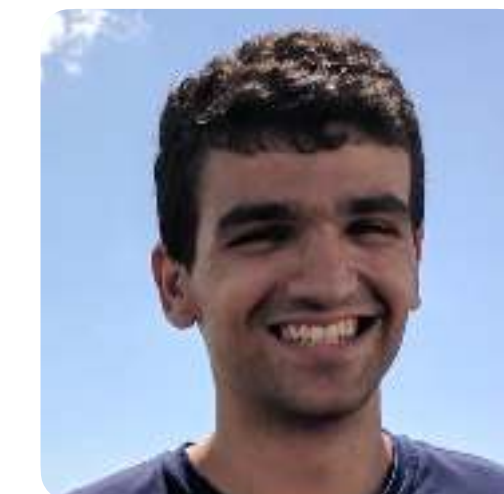
data



$$D = \{x^i\}_{i=1}^N$$



how can we bound
the true probability?



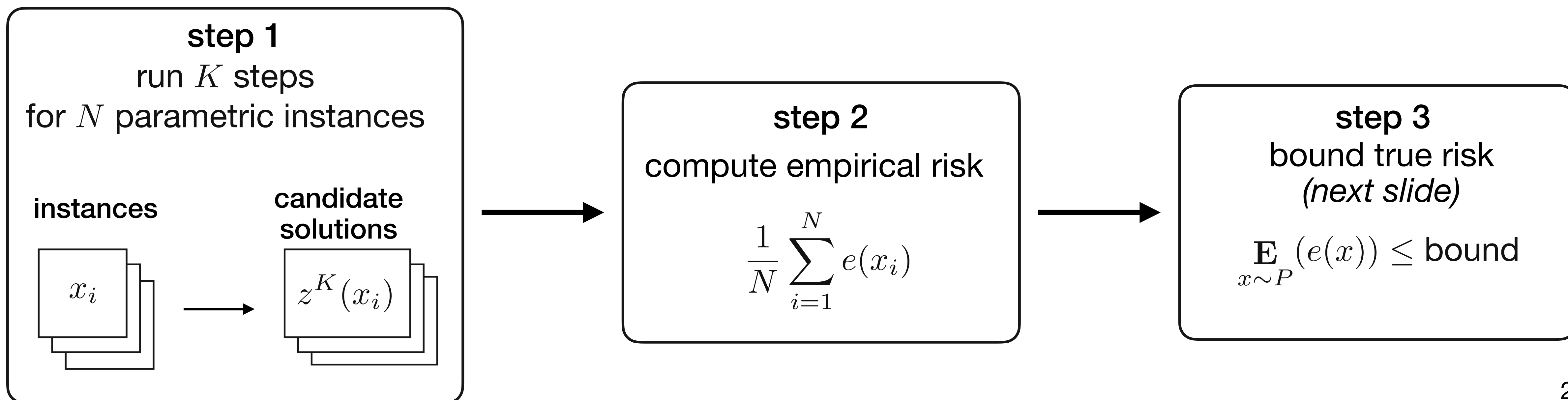
Our recipe to bound performance

goal

*estimate probability of
computing bad-quality solutions*

$$\mathbf{P}(r^K(x) > \epsilon) = \mathbf{E}_{x \sim P}(e(x))$$

\uparrow any metric (e.g., fixed-point residual) \uparrow error $\mathbf{1}(r^K(x) > \epsilon)$



Statistical learning gives us probabilistic guarantees

sample convergence bound
with probability $1 - \delta$

$$\mathbf{P}(r^K(x) > \epsilon) = \mathbf{E}_{x \sim P}(e(x)) \leq \text{kl}^{-1} \left(\frac{1}{N} \sum_{i=1}^N e(x_i) \mid \frac{\log(2/\delta)}{N} \right)$$

true risk
inverse kl divergence
(1D convex problem)
empirical risk
regularizer
number of instances

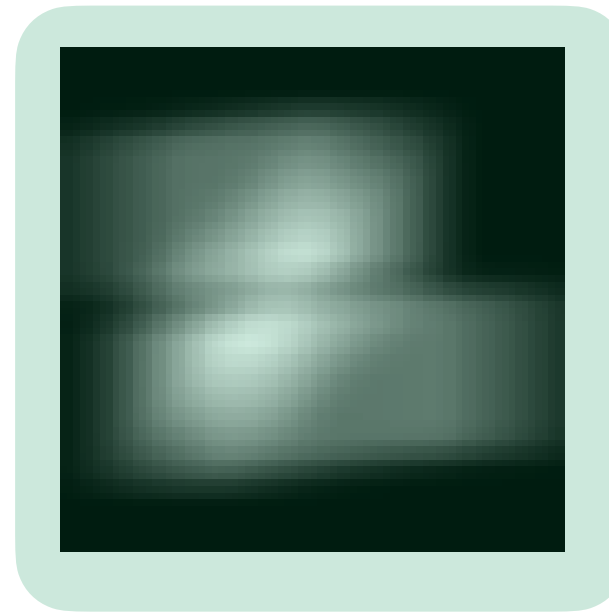
interpretation of bound equal to B

With probability $1 - \delta$, the fixed-point residual is above ϵ after K steps

B fraction of times

Success rates for OSQP in image deblurring

minimize $\|Az - x\|_2^2 + \lambda \|z\|_1$
 subject to $0 \leq z \leq 1$

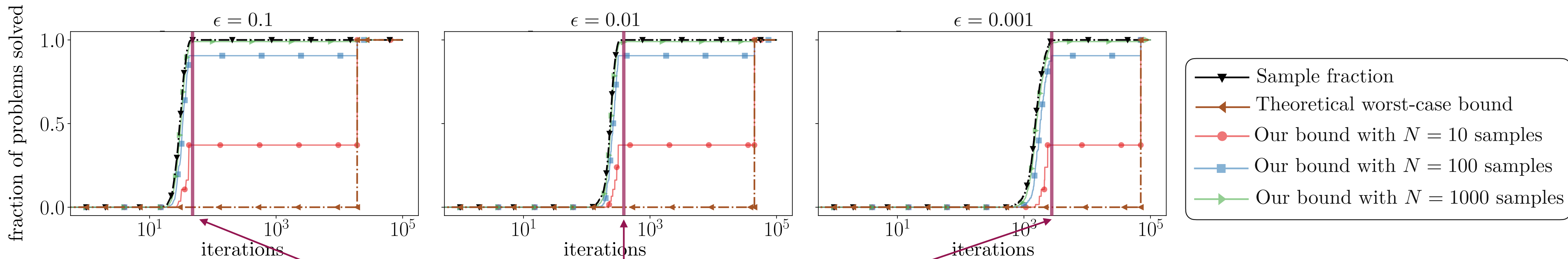


deblurred image

blurred image

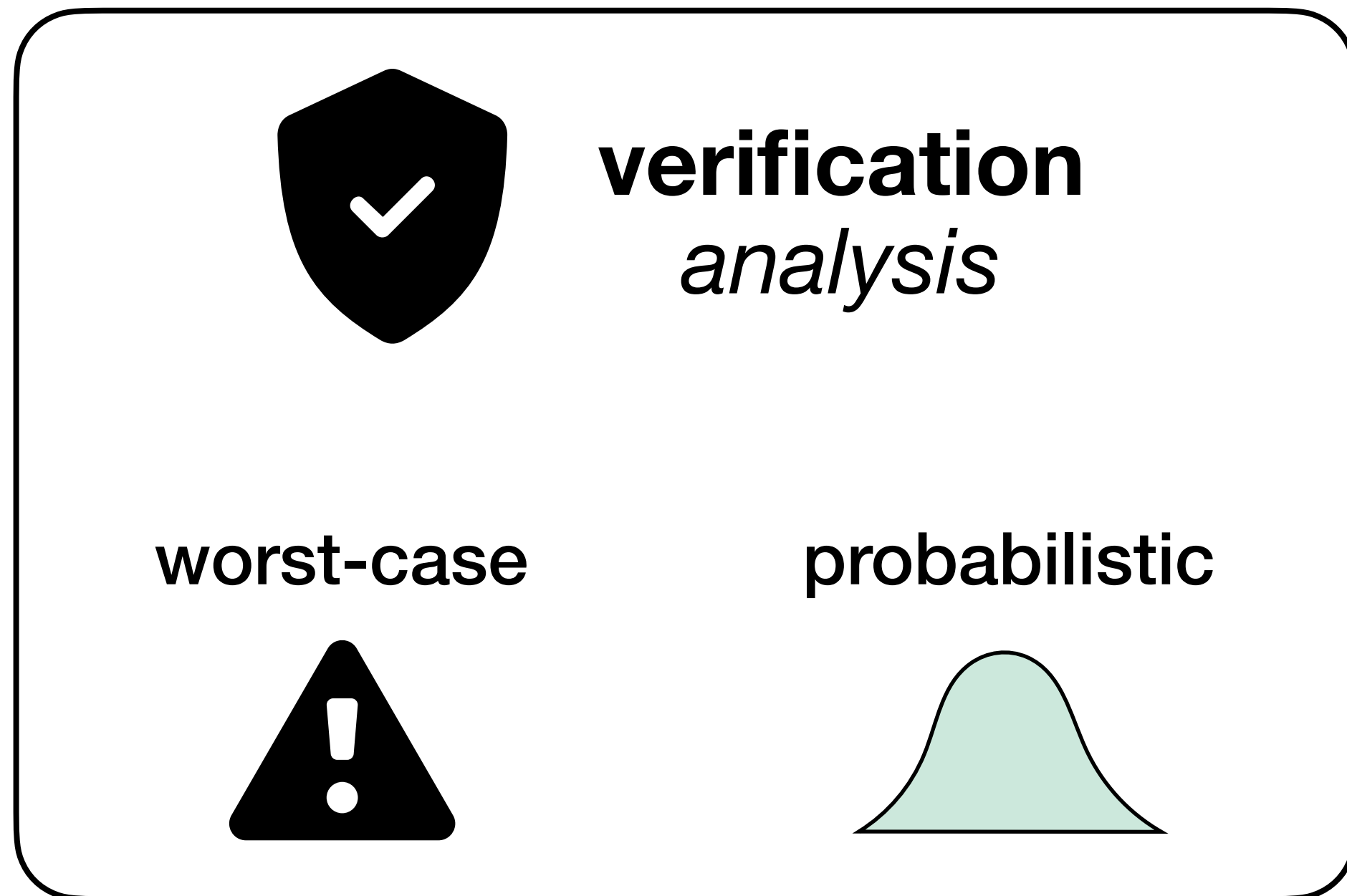
fraction of problems solved

$$1 - \mathbf{E}_{x \sim P} (e(x)) \leftarrow \mathbf{1}(r^K(x) > \epsilon)$$



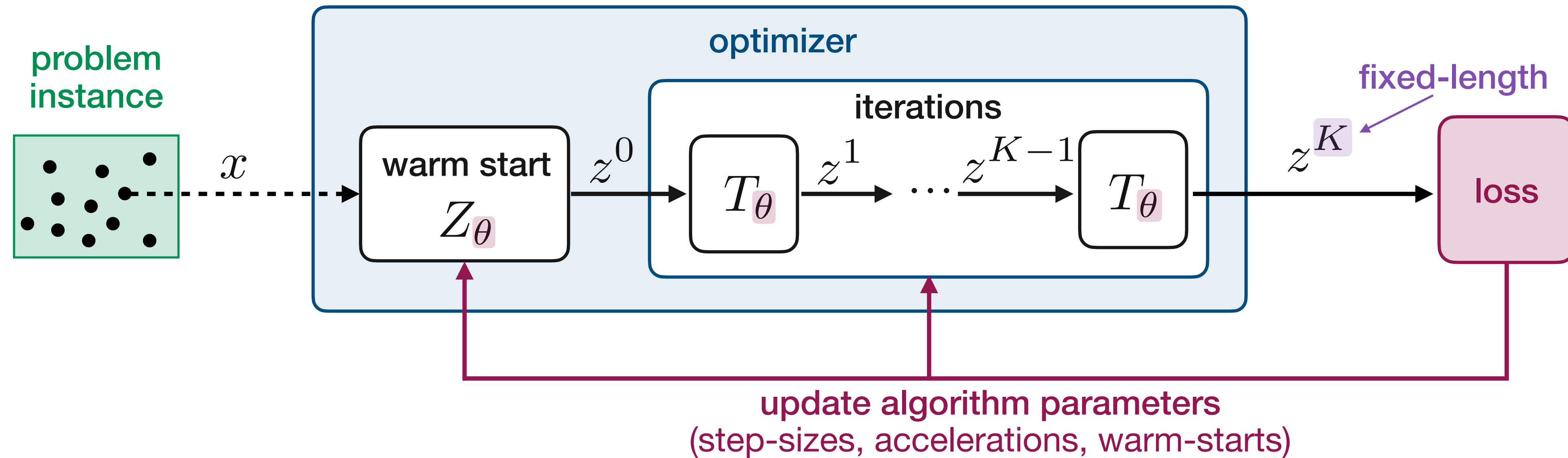
iterations required to solve all test instances

First-order methods in parametric convex optimization



Algorithm design

Training algorithms as fixed-length computational graphs



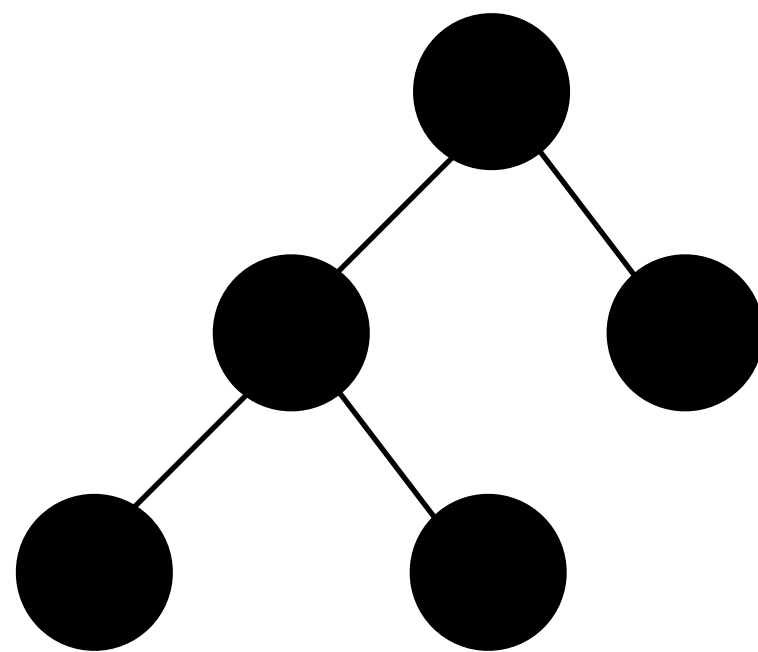
example
projected gradient descent

$$z^{k+1} = \Pi_{C(x)}(z^k - \theta \nabla_z f(z^k, x))$$

Learning can accelerate optimizers

Combinatorial optimization

B. Dilkina, E. Khalil, A. Lodi, P. Van Hentenryck, P. Bonami, S. Jegelka, ...



our previous contributions

The voice of optimization

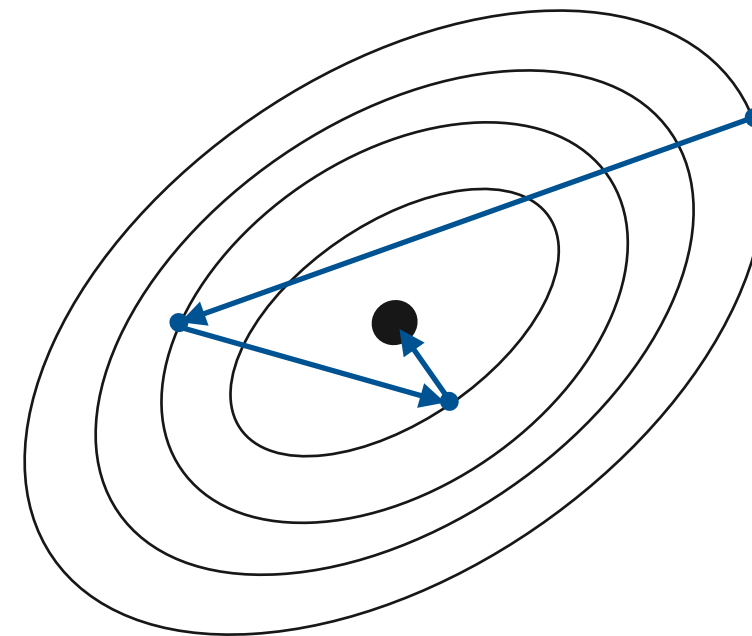
D. Bertsimas, B. Stellato
Machine Learning (2021)

Online mixed-integer optimization in milliseconds

D. Bertsimas, B. Stellato
INFORMS Journal on Computing (2022)

Continuous optimization

W. Yin, B. Amos, Z. Kolter, M. Andrychowicz, C. Finn, P. Van Hentenryck ...



our previous contributions

Accelerating quadratic optimization with reinforcement learning

J. Ichnowski, P. Jain, B. Stellato, ... et al.
NeurIPS (2021)

No performance guarantees



Can we build rigorous and data-driven performance guarantees?

Statistical learning theory for optimization algorithms

supervised learning

learning to optimize

input



problem instance
(with parameter x)

hypothesis

cat

residual $r_{\theta}^K(x)$

error

0 (1 if wrong)

$$e_{\theta}(x) = \mathbf{1}(r_{\theta}^K(x) > \epsilon)$$

guarantees

expected loss
on **new data**

expected loss
on **new problem instances**

algorithm parameters
(step-sizes,
accelerations,
warm-starts)

PAC-Bayes generalization bounds

learning task

$$\text{minimize}_{\Theta} \mathbf{E}_{\theta \sim \Theta} \mathbf{E}_{x \sim P} (e_{\theta}(x))$$

distribution of algorithm parameters
(step-sizes, accelerations, warm-starts)

can be anything

1. Pick *prior* Θ_0 before observing data

2. Observe data $D = \{x^i\}_{i=1}^N$

3. Learn *posterior* Θ : $\theta \sim \Theta$

4. Bound performance $\mathbf{P}^N \left(\mathbf{E}_{\theta \sim \Theta} \mathbf{E}_{x \sim P} (e_{\theta}(x)) \leq \hat{t}_N \right) \geq 1 - \delta$

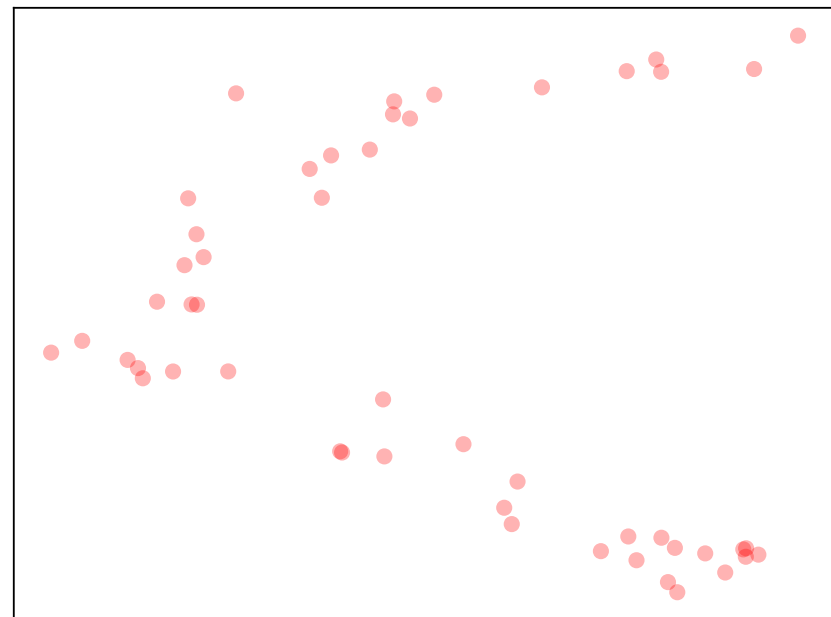
McAllester (1999), Maurer (2004)

data-driven bound

$$\hat{t}_N = \text{kl}^{-1} \left(\underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{E}_{\theta \sim \Theta} (e_{\theta}(x_i))}_{\text{empirical risk}} \mid \underbrace{\frac{\text{KL}(\Theta \parallel \Theta_0) + \log(2\sqrt{N}/\delta)}{2N}}_{\text{regularizer}} \right)$$

Robust Kalman Filtering with learned warm starts

noisy trajectory



$$x = \{y_t\}_{t=0}^{T-1}$$



second-order cone program solver (SCS)

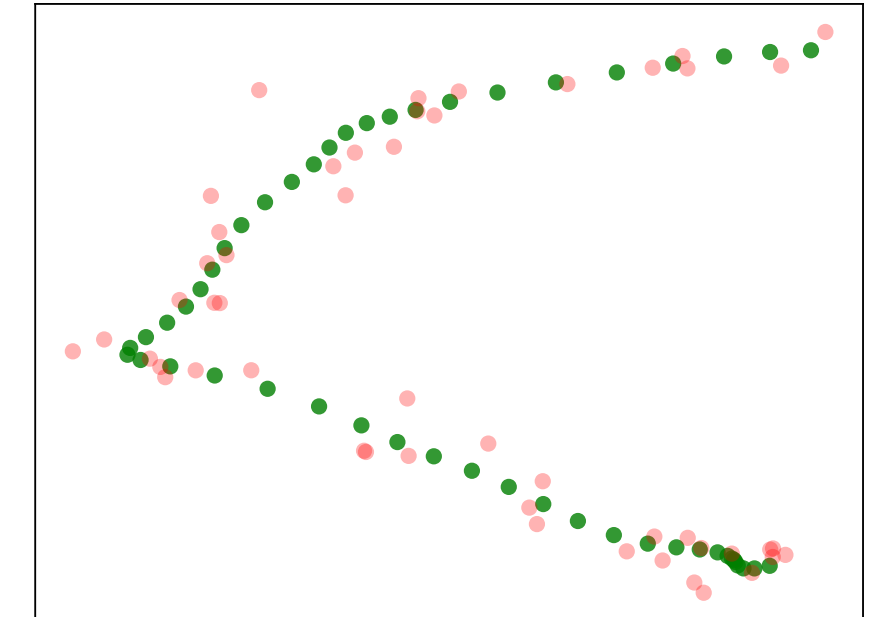
minimize $\sum_{t=0}^T \|w_t\|_2^2 + \psi(v_t)$ Huber loss

subject to $s_{t+1} = As_t + Bw_t, \quad t = 0, \dots, T - 1$

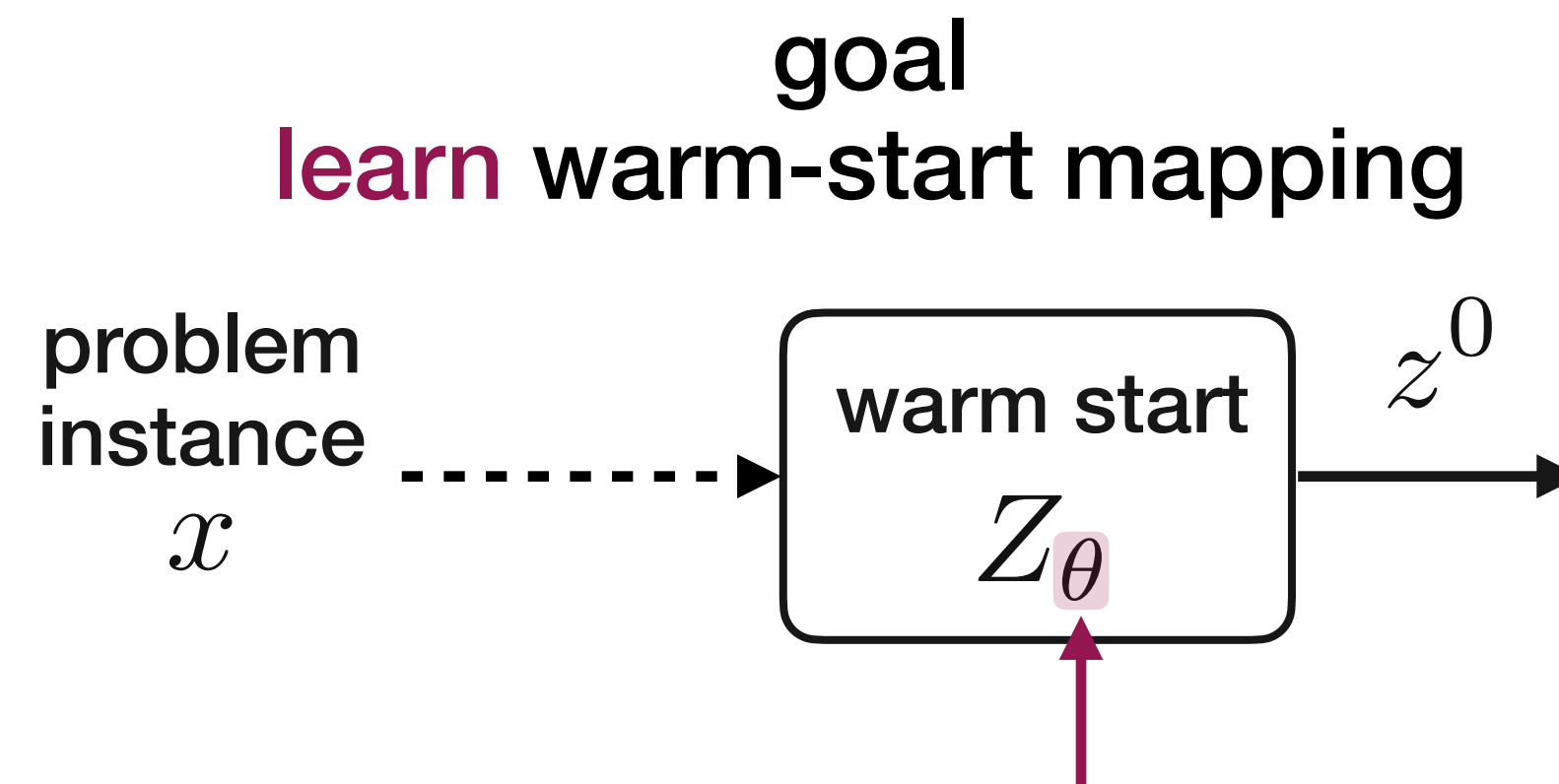
$y_t = Cs_t + v_t, \quad t = 0, \dots, T$



recovered trajectory

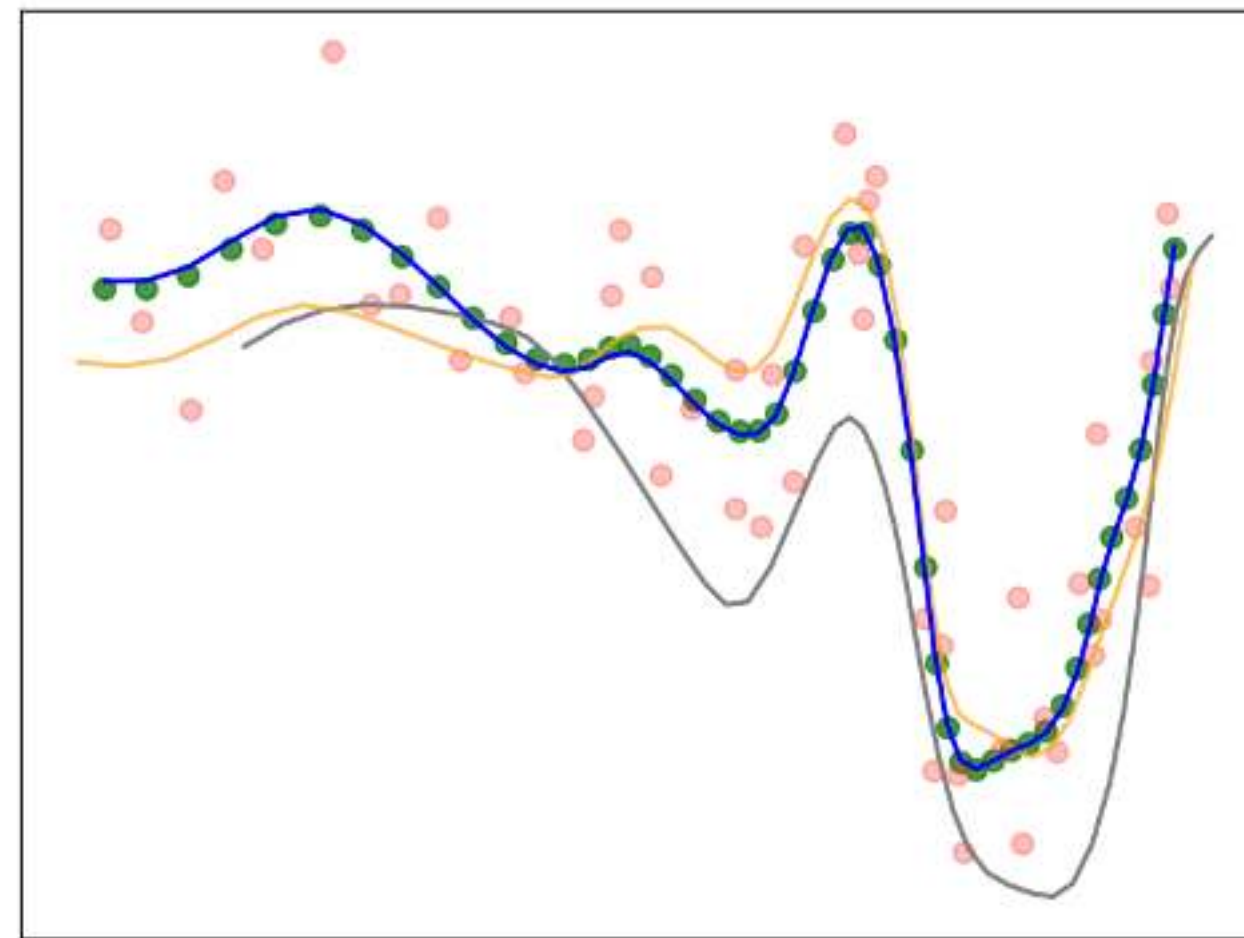
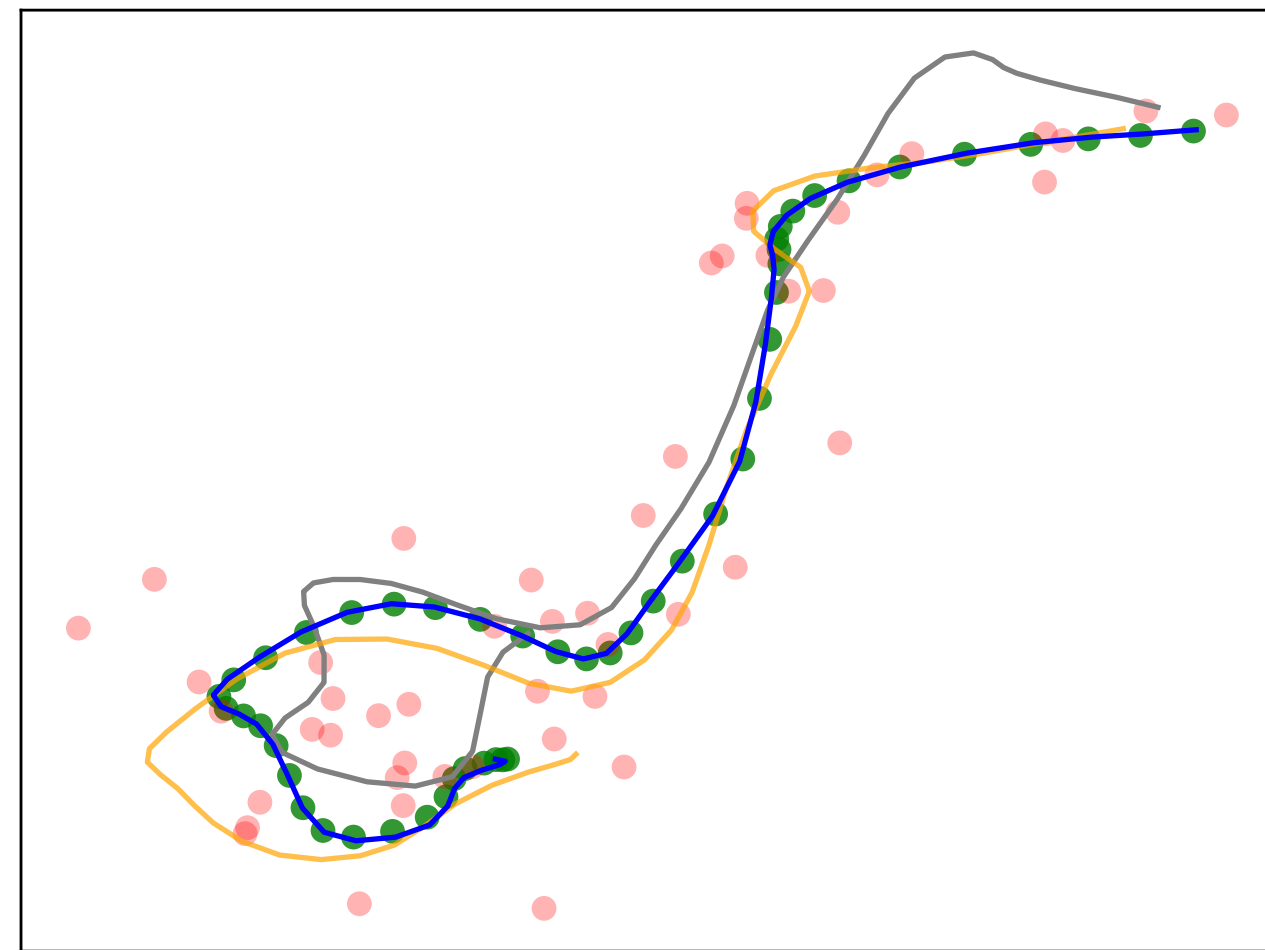


$$z^* = \{s_t^*, w_t^*, v_t^*\}_{t=0}^{T-1}$$



Robust Kalman Filtering with learned warm starts

two example trajectories



points

- noisy trajectory
- optimal solution

Solution after 5 fixed-point iterations
with different warm-starts

- nearest neighbor
- previous solution
- learned $K = 5$

with learning, we can
estimate the state well

we also showed
warm-start specific
PAC Bayes generalization
guarantees



Learning to Warm-Start Fixed-Point Optimization Algorithms

R. Sambharya, G. Hall, B. Amos, and B. Stellato

Journal of Machine Learning Research (2024)

github.com/stellatogrp/l2ws

Signal reconstruction with learned optimizer

minimize $\|Dz - x\|_2^2 + \lambda \|z\|_1$

known dictionary (pointing to D)

noisy signal (pointing to x)

reconstructed signal (pointing to z)

performance metric
normalized mean squared error

$$\text{NMSE}_{\text{dB}}(z) = 10 \log_{10} \left(\|z - \bar{z}\|^2 / \|\bar{z}\|^2 \right)$$

ground truth (pointing to \bar{z})

classical algorithm (ISTA)

$$z^{k+1} = \phi_{\lambda t} \left(z^k - t2D^T (Dz^k - x) \right)$$

shrinkage operator

$$\phi_{\lambda t}(v) = \max\{v, \lambda t\} - \max\{-v, \lambda t\}$$

learned variants (e.g., ALISTA)

$$z^{k+1} = \phi_{\gamma^k} \left(z^k - \psi^k W^T (Dz^k - x) \right)$$

algorithm parameters

$$\theta = \left\{ \gamma^k, \psi^k \right\}_{k=0}^{K-1}$$

Conclusions

Algorithm Design and Verification for Parametric Convex Optimization


1. **parametric** structure matters

2. **data** can help us

- **design** optimization algorithms 
- **verify** their performance 

3. we should **rethink optimization algorithms**


traditional view



- general-purpose
- one-size-fits all



new view



- task-specific
- trainable
- deployable anywhere

Backup

PAC-Bayes generalization guarantees for learned warm starts

β -contractive case

$$\|Tx - Ty\|_2 \leq \beta \|x - y\|_2 \quad \forall x, y$$

$$\beta \in (0, 1)$$

Theorem: for any $\gamma > 0$ with probability at least $1 - \delta$

$$\mathbf{E}_{x \sim \mathcal{X}} \ell_{\theta}^k(x) \leq \frac{1}{N} \sum_{i=1}^N \ell_{\theta}^k(x_i) + 2\beta^k \gamma + \mathcal{O} \left(\frac{\beta^k}{\gamma} (2D + 1) \sqrt{\frac{c_2(\theta) + \log(\frac{LN}{\delta})}{N}} \right)$$

risk

empirical risk

penalty term

bound on $\|z^*(x)\|_2$

As the number of iterations $k \rightarrow \infty$ the penalty term goes to zero

The contractive factor β directly affects the penalty term

We combine operator theory with PAC-Bayes theory to get the bound

Computing the KL Inverse with Convex Optimization

KL divergence between Bernoulli distributions
 $\text{kl}(q \parallel p) := \text{KL}(\text{Bernoulli}(q), \text{Bernoulli}(p))$

Many PAC-Bayes-type bounds bound the **risk** implicitly

$$\text{kl}(q \parallel p) \leq c$$

empirical risk risk regularizer

Inverting the KL divergence

$$p^* = \text{kl}^{-1}(q \mid c) = \begin{array}{l} \text{maximize } p \\ \text{subject to } q \log\left(\frac{q}{p}\right) + (1 - q) \log\left(\frac{1-q}{1-p}\right) \leq c \\ 0 \leq p \leq 1 \end{array}$$

Solved within a
millisecond on my laptop