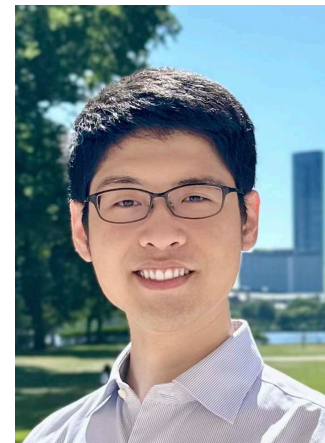


# Level-Set Geometry and the Performance of Restarted-PDHG for Conic LP

Zikai Xiong (with Robert Freund)

Optimization Workshop: Theory, Algorithms and Application

Universidad de los Andes, Bogota



Zikai Xiong  
(MIT OR Center)



Robert Freund  
(MIT Sloan)

Bogota, Colombia

December 2024

# Huge-scale optimization is everywhere

**Manufacturing**



**Machine Learning**



**Energy**



**Transportation**



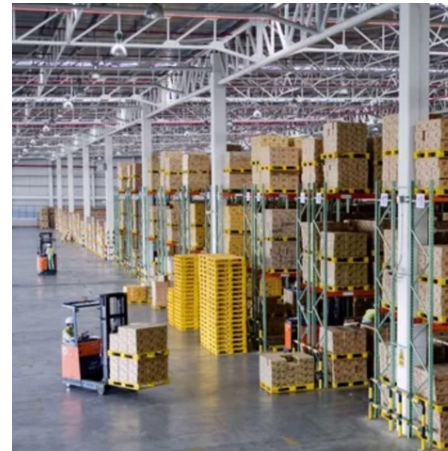
**Healthcare**



**Markets and Auctions**



**Supply Chains**



**Agriculture**

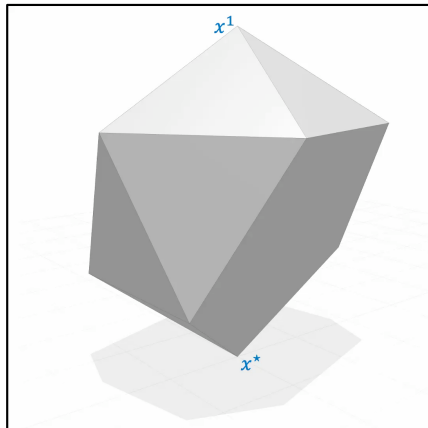


# History of Linear Optimization (“LO” or “LP”)

1947

**Simplex  
Method**

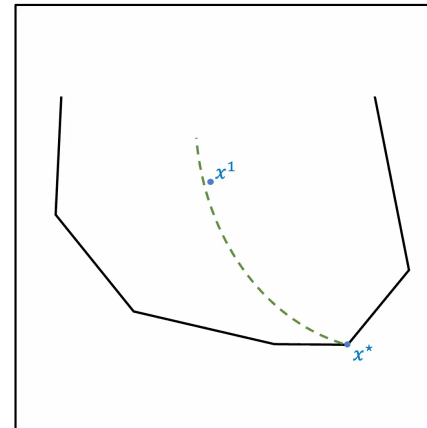
[George Dantzig, 1947]  
75+ years ago



1984

**Interior Point  
Method**

[Narendra Karmarkar, 1984]  
40 years ago



# Simplex and IPMs require expensive matrix factorizations

Consider an LP instance with  $n$  decision variables and  $\frac{n}{2}$  linear constraints, whose constraint matrix has sparsity = 0.05

**Number of  
variables ( $n$ )**

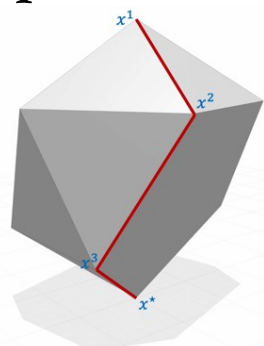
**Cost of one IPM  
iteration**

Hence the emergence of FOMs for solving huge (and also not-so-huge) LP instances

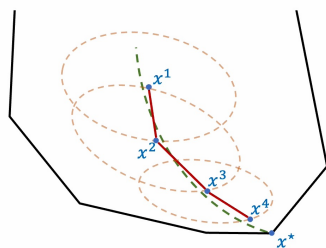
# Recent Advances on Huge-Scale LP Solvers in the Industry

## Classic methods

### Simplex method



### Interior-point method



## First-order methods

- Primal-Dual Hybrid Gradient (“PDHG”, “Chambolle-Pock method”)
- Tackles huge-scale problems
- Benefits from modern computational architecture (such as GPU)



2021



February, 2024



March, 2024



March, 2024



April, 2024



October, 2024

**We are witnessing a dramatic shift from classic methods to first-order methods**

# Huge-Scale LP Research

- SCS: Operator splitting/ADMM [O'Donoghue, Chu, Parikh, Boyd, 2016]
- ABIP+: ADMM-based interior-point method [Lin, Ma, Ye, Zhang, 2021] & [Deng, et al., 2022]
- Semi-smooth Newton augmented Lagrangian [Li, Sun, Toh, 2020]
- **Primal-Dual Hybrid Gradient (PDHG)** with restarts, applied directly to the primal-dual saddle point problem [Applegate, Hinder, Lu, Lubin, 2023] & [Applegate, et al., 2021] (**2024 Beale-Orchard-Hays Prize**)
- **GPU implementations** of PDHG in Julia and C [Lu and Yang, 2023] & [Lu et al., 2023]
- **Guarantees for PDHG for LP** using “Limiting Error Ratios” and LP Sharpness [Xiong and F 2023]
- **Guarantees for PDHG for CLP** – using level-set geometry [Xiong and F 2024]

# This talk is based on material from three papers

## **For LP and conic optimization:**

- New computational guarantees based on problem (sub)level-set geometry

Xiong, Z., and Freund, R. M. (2024). The Role of Level-Set Geometry on the Performance of PDHG for Conic Linear Optimization.

## **For LP with unique optima:**

- Closed-form iteration bound
- Two-stage performance of PDHG
- “Average-case” polynomial-time complexity guarantee

Xiong, Z. (2024). Accessible Theoretical Complexity of the Restarted Primal-Dual Hybrid Gradient Method for Linear Programs with Unique Optima.

Xiong, Z. (2024). Probabilistic Analysis of Restarted PDHG for Linear Programming Problems (working paper).

# Conic Linear Optimization (“CLO” or “CLP”)

## CLP in standard form

(primal)

$$\begin{aligned} \min \quad & c^\top x \\ \text{s.t.} \quad & Ax = b \\ & x \in \mathcal{K} \end{aligned}$$

(dual)

$$\begin{aligned} \max \quad & b^\top y \\ \text{s.t.} \quad & c - A^\top y \in \mathcal{K}^* \end{aligned}$$

Decision variables

- $x \in R^n$  (for primal problem)
- $y \in R^m$  (for dual problem)

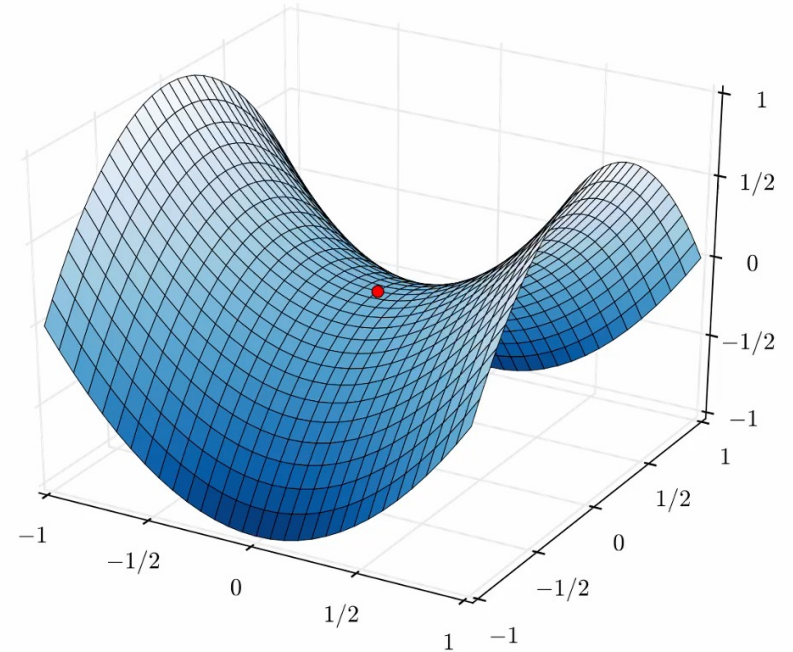
## CLP saddlepoint formulation

$$\min_{x \in \mathcal{K}} \max_y c^\top x - y^\top Ax + b^\top y$$

# Primal-Dual Hybrid Gradient Method (PDHG)

Conic Optimization in  
Saddlepoint Form

$$\min_{x \in \mathcal{K}} \max_y c^\top x + b^\top y - y^\top A x$$



PDHG

$$x^{k+1} \leftarrow \text{Proj}_{\mathcal{K}} \left( x^k - \tau (c - A^\top y^k) \right)$$

Gradient w.r.t.  $x^k$

$$y^{k+1} \leftarrow y^k + \sigma (b - A x^{k+1}) - \sigma A (x^{k+1} - x^k)$$

Gradient w.r.t.  $y^k$

Momentum Term

# Primal-Dual Hybrid Gradient for Conic Optimization

## PDHG

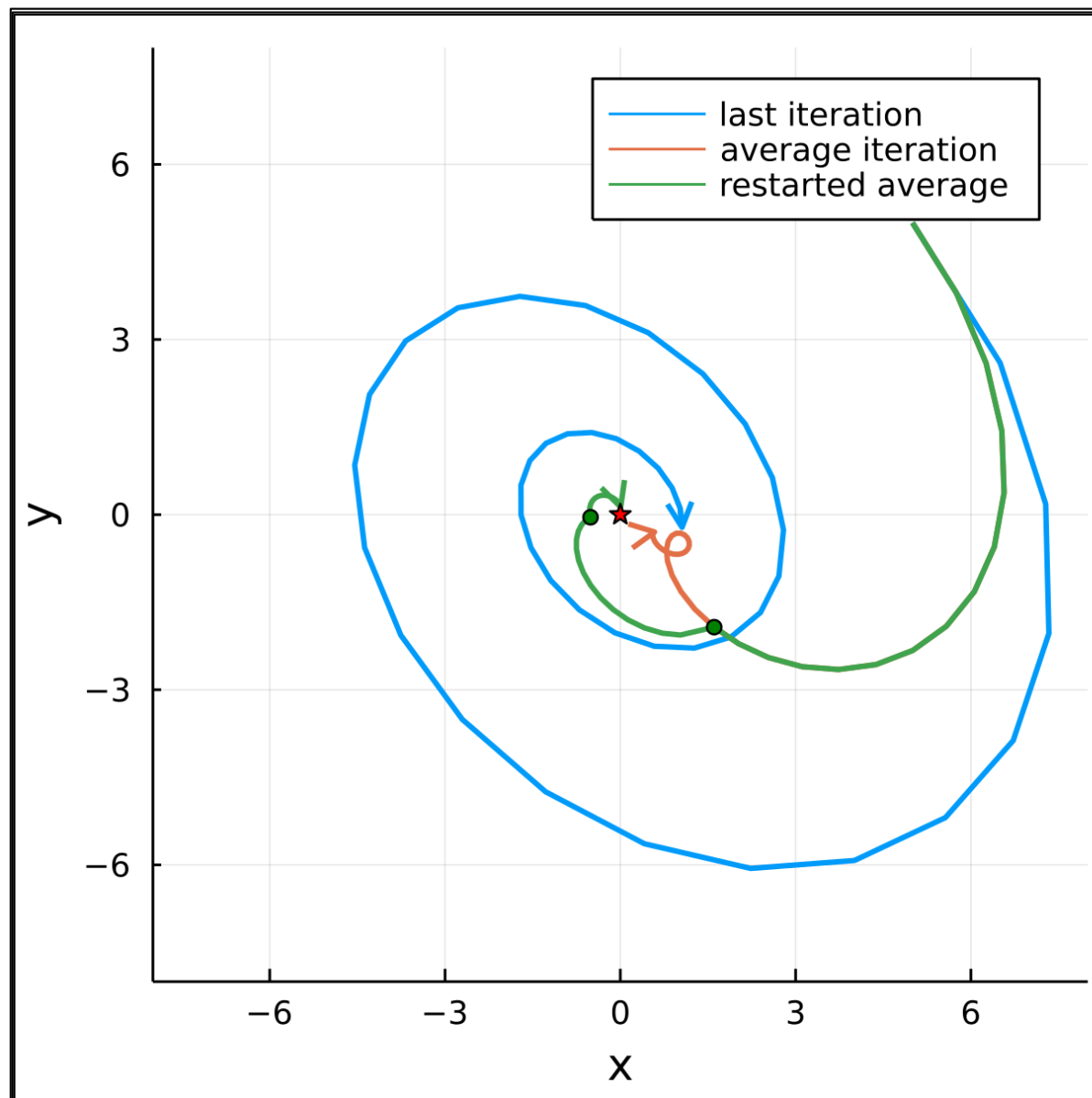
$$\begin{aligned}x^{k+1} &\leftarrow \text{Proj}_{\mathcal{K}} \left( x^k - \tau(c - A^\top y^k) \right) \\y^{k+1} &\leftarrow y^k + \sigma(b - Ax^{k+1}) - \sigma A(x^{k+1} - x^k)\end{aligned}$$

- **Inexpensive iterations:**  
Only requires matrix-vector multiplications
- **“Fast” convergence rates:**  
Adaptive restarts based on average iterates yield global linear convergence on LP [Applegate, Hinder, Lu, Lubin, 2023]

**We use “PDHG” to denote “PDHG with adaptive restarts”**

# Motivation for Restarts for PDHG: “Visualization”

$$\min_x \max_y x \cdot y$$

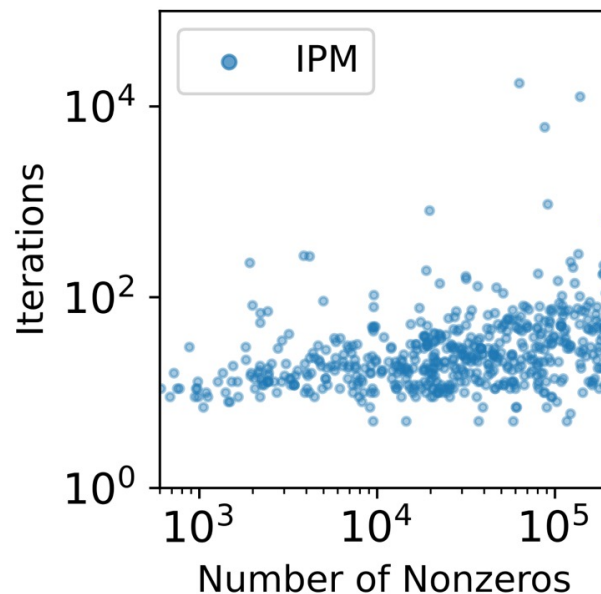


\*figure courtesy Haihao Lu

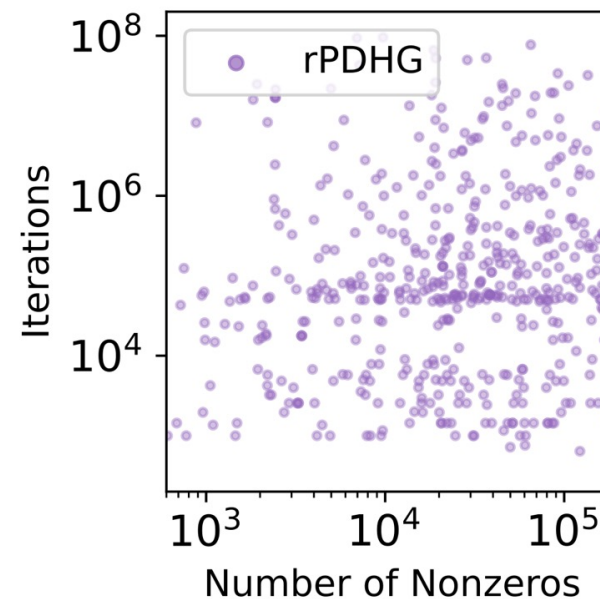
# Challenge I: Variability in the Performance of PDHG

- PDHG uses many more iterations than an IPM  
makes sense, it is a first-order method ... IPM iterations are hugely expensive while PDHG iterations are very cheap
- Some small problem instances require a very large number of PDHG iterations  
a real challenge for PDHG

IPM Iterations needed for  
LP relaxations from MIPLIB 2017



PDHG iterations  
LP relaxations from MIPLIB 2017



# A seemingly easy LP instance

For  $\gamma > 0$ , define:

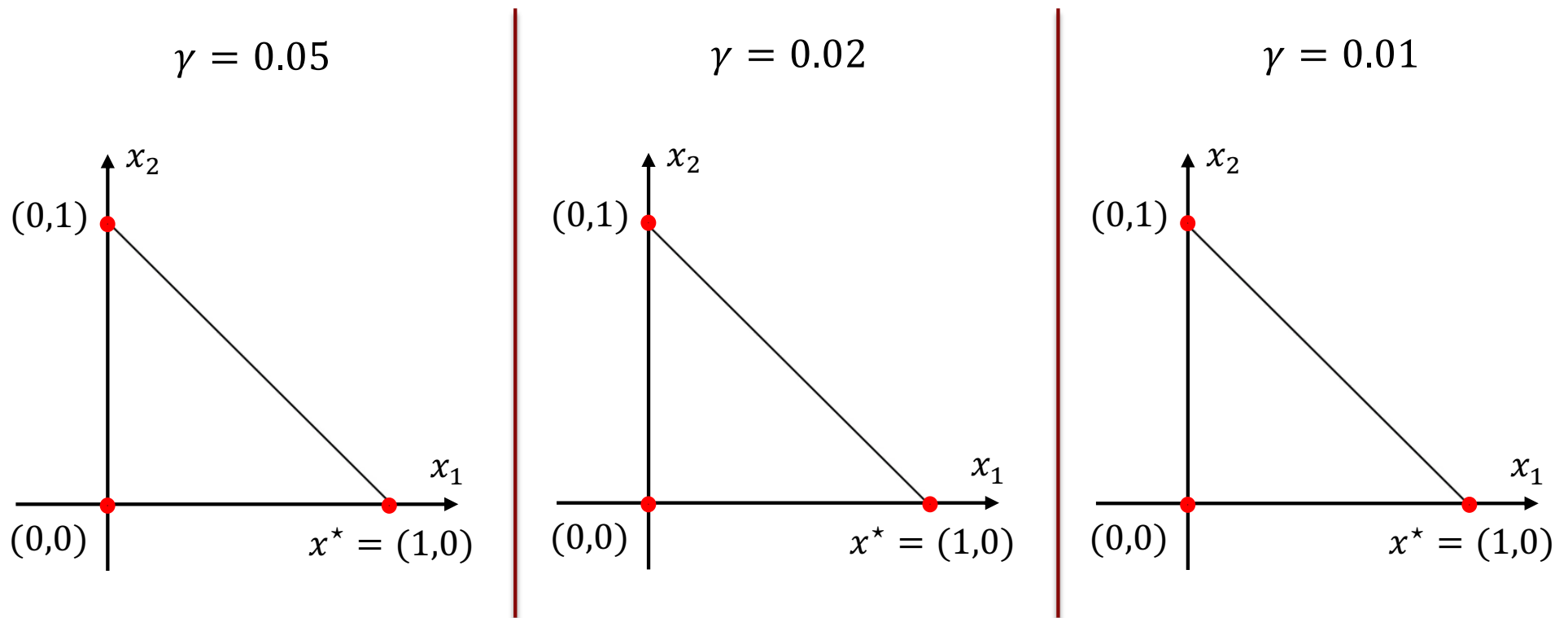
$$P(\gamma): \quad \begin{array}{ll} \min_{x_1, x_2} & -(1 + \gamma)x_1 - x_2 \\ \text{s. t.} & x_1 + x_2 = 1 \\ & x_1 \geq 0, x_2 \geq 0 \end{array}$$

For  $P(\gamma)$ , the optimal solution is always  $x^* = (1, 0)$

$P(\gamma)$  is always easy for the simplex method and interior-point methods

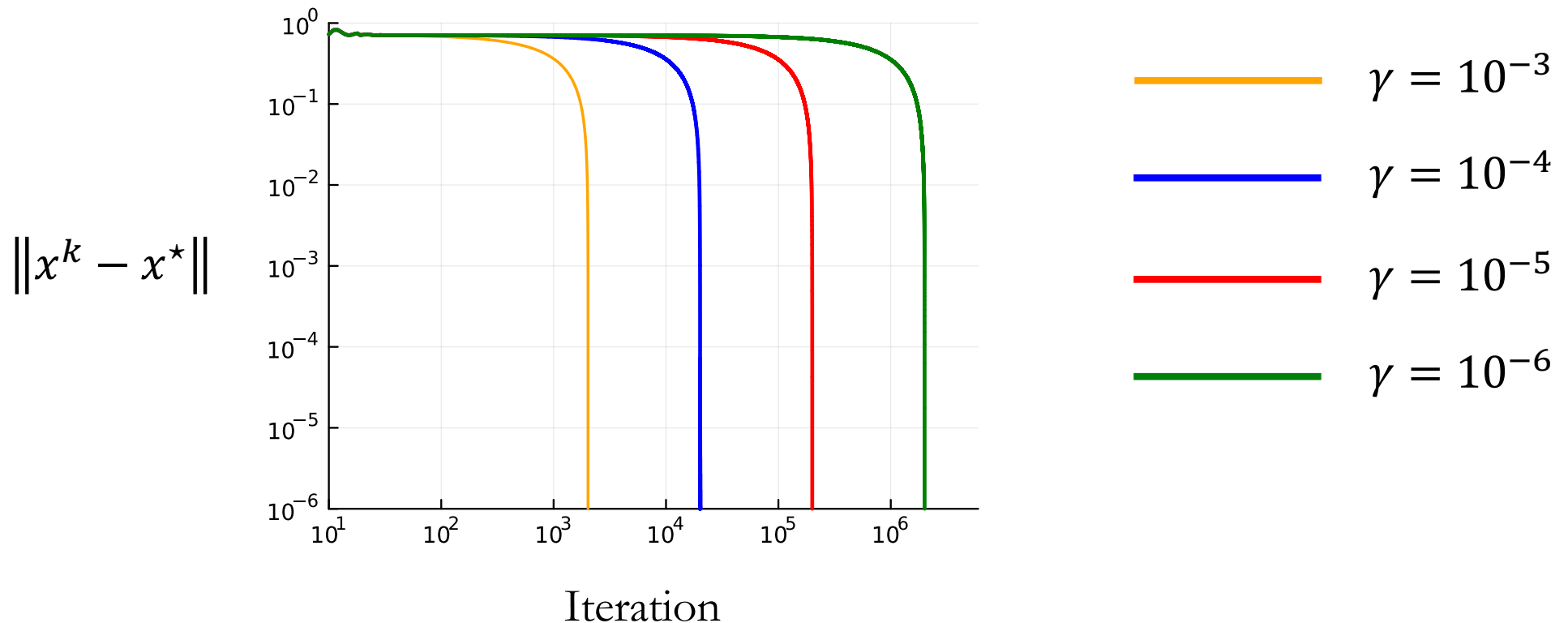
# A seemingly easy LP instance

The trajectory of PDHG is as follows



1. The trajectory of PDHG is of course quite different from simplex or IPM.
2. The smaller the  $\gamma$  is, the slower the convergence will be.

# A seemingly easy LP instance



When  $\gamma = 10^{-6}$ , PDHG requires more than 1,000,000 iterations.  
What **conditions** of  $P(\gamma)$  make it so hard for PDHG?

# Challenge II: Loose/unworkable computational guarantees

Existing computational guarantees:

**Theorem** [Applegate, Hinder, Lu, Lubin, 2023] PDHG computes an  $\varepsilon$ -optimal solution within

$$O\left((\|x^*\| + \|y^*\|) \cdot \|A\| \cdot H(K) \cdot \log\left(\frac{\|x^*\| + \|y^*\|}{\varepsilon}\right)\right)$$

iterations.

$H(K)$  is the global Hoffman constant of the matrix  $K$  of the KKT system

$$K = \begin{pmatrix} A & & & \\ -A & & & \\ & A^\top & & \\ c^\top & & -b^\top & \end{pmatrix}, \quad H(K) \approx \frac{1}{\min_{J \subseteq \{1, \dots, 2m+n+1\}} \sigma_{\min}^+(K_J)}$$

$m$ : the number of constraints  
 $n$ : the number of variables

1. Usually too loose
2. Very hard to compute/validate/analyze

## Challenge II: Loose/unworkable computational guarantees

**Theorem** [Applegate, Hinder, Lu, Lubin, 2023] PDHG computes an  $\varepsilon$ -optimal solution within

$$O\left(\left(\|x^*\| + \|y^*\|\right) \cdot \|A\| \cdot \mathbf{H}(\mathbf{K}) \cdot \log\left(\frac{\|x^*\| + \|y^*\|}{\varepsilon}\right)\right)$$

iterations.

Key questions:

- What conditions of the problem actually drive the performance of PDHG? **Sublevel-set Geometry**
- Can we improve these conditions and so improve computational performance in theory/practice?  
**Yes, by using Hessian rescaling (not covered in this talk)**
- Can we shrink the gap between theory and practice?  
**Yes, by using probabilistic average case analysis (if time permits...)**

# Challenge II: Loose/unworkable computational guarantees

Existing computational guarantees:

**Theorem** [Applegate, Hinder, Lu, Lubin, 2023] PDHG computes an  $\varepsilon$ -optimal solution within

$$O\left((\|x^*\| + \|y^*\|) \cdot \|A\| \cdot H(K) \cdot \log\left(\frac{\|x^*\| + \|y^*\|}{\varepsilon}\right)\right)$$

iterations.

$H(K)$  is the global Hoffman constant of the matrix  $K$  of the KKT system

$$K = \begin{pmatrix} A & & & \\ -A & & & \\ & A^\top & & \\ c^\top & & -b^\top & \end{pmatrix} \quad H(K) \approx \frac{1}{\min_{J \subseteq \{1, \dots, 2m+n+1\}} \sigma_{\min}^+(K_J)}$$

Likely way too loose, and very hard to compute/validate/analyze.

# Challenge II: Loose/unworkable computational guarantees

Existing computational guarantees:

**Theorem** [Applegate, Hinder, Lu, Lubin, 2023] PDHG computes an  $\varepsilon$ -optimal solution within

$$O\left(\left(\|x^*\| + \|y^*\|\right) \cdot \|A\| \cdot H(K) \cdot \log\left(\frac{\|x^*\| + \|y^*\|}{\varepsilon}\right)\right)$$

iterations.

Key question:

- What conditions of the problem actually drive the performance of PDHG?

**Sublevel-set Geometry**

# Sublevel-set geometry and new performance guarantees for PDHG

# Primal-Dual Slack Space

Primal

$$\min_{\boldsymbol{x}} c^{\top} \boldsymbol{x}$$

$$\text{s. t. } A\boldsymbol{x} = \boldsymbol{b}$$

$$\boldsymbol{x} \in \mathcal{K}$$

Dual

$$\max_{\boldsymbol{y}, \boldsymbol{s}} b^{\top} \boldsymbol{y}$$

$$\text{s. t. } A^{\top} \boldsymbol{y} + \boldsymbol{s} = \boldsymbol{c}$$

$$\boldsymbol{s} \in \mathcal{K}^*$$

The “primal-dual slack-space variable” is  $\boldsymbol{w}$  :

$\boldsymbol{w} := (\boldsymbol{x}, \boldsymbol{s})$  are primal/dual feasible slacks

Duality gap:  $\text{Gap}(\boldsymbol{x}, \boldsymbol{s}) = c^{\top} \boldsymbol{x} - b^{\top} \boldsymbol{y}$

(which is a linear function of  $\boldsymbol{x}$  and  $\boldsymbol{s}$ )

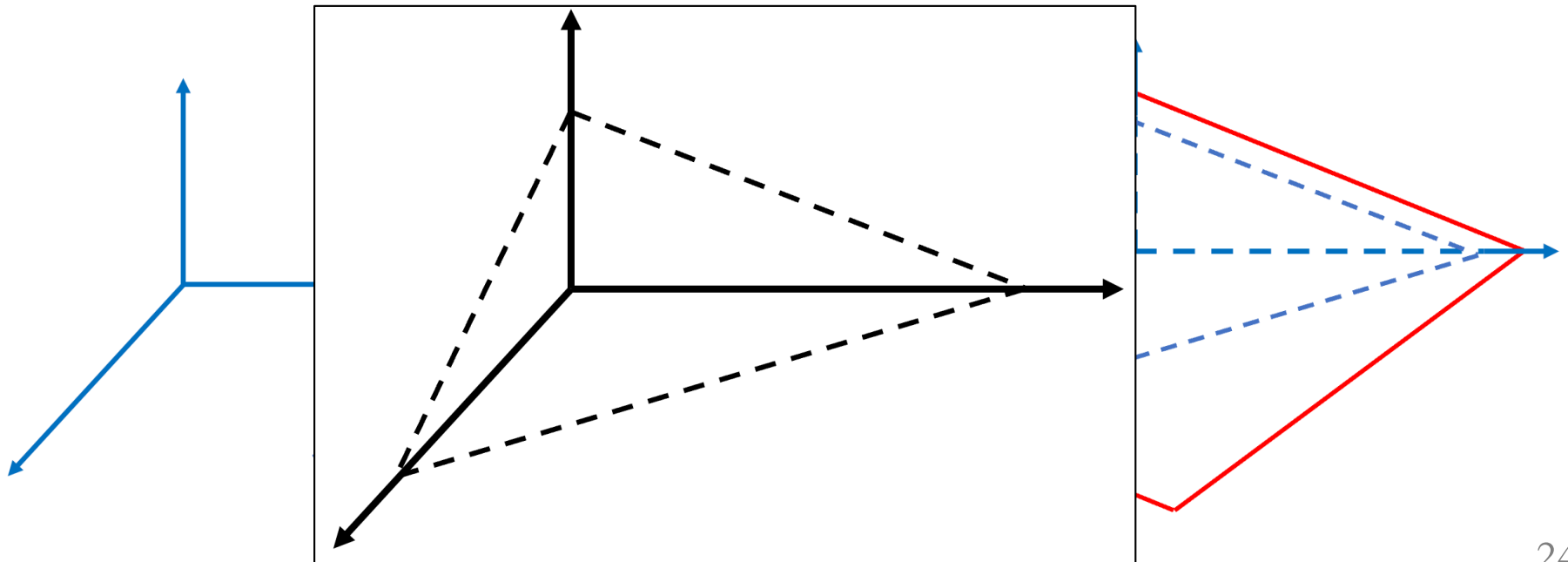
# The feasible primal-dual slack-space variables

$$\begin{aligned} \min_{x} \quad & c^T x \\ \text{s. t.} \quad & Ax = b \\ & x \in \mathcal{K} \end{aligned}$$

$$\begin{aligned} \max_{y,s} \quad & b^T y \\ \text{s. t.} \quad & A^T y + s = c \\ & s \in \mathcal{K}^* \end{aligned}$$

$(x, s)$  in the primal and dual cone for CLP

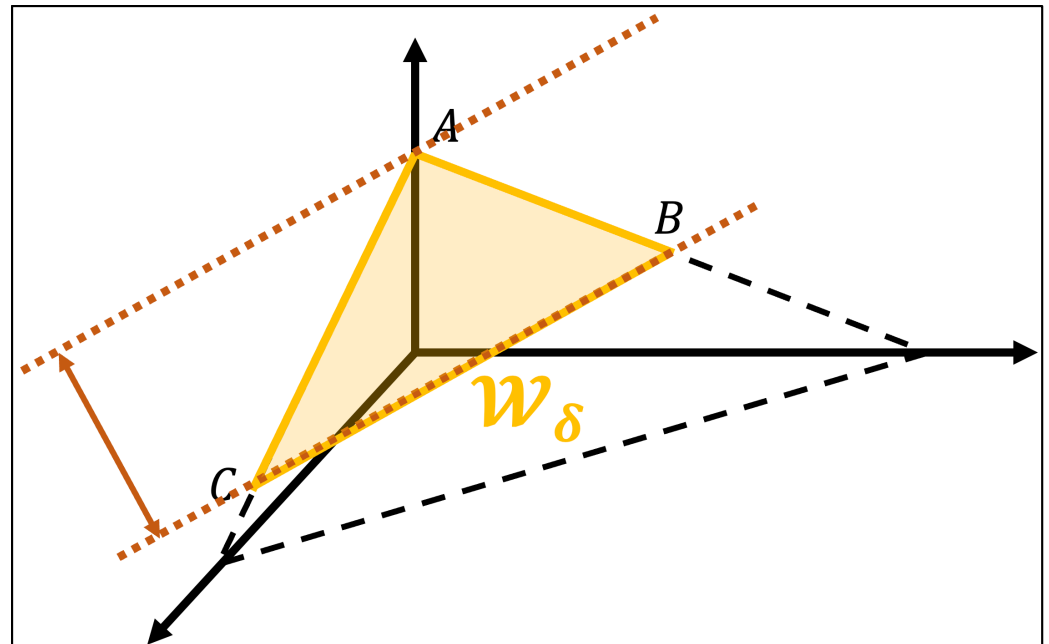
$(x, s)$  lies in an affine subspace



# Primal-Dual Slack Sublevel Set

$$\mathcal{W}_\delta := \left\{ w := (x, s) \mid \begin{array}{l} w \text{ is primal/dual feasible} \\ \mathbf{Gap}(w) \leq \delta \end{array} \right\}$$

Note:  $\mathcal{W}_0 = \mathcal{W}^*$



# Worst-case complexity of PDHG (under unique optima)

**Theorem** [Xiong and F 2024]: Suppose  $w^*$  is unique. PDHG computes an  $\varepsilon$ -optimal solution within

$$\tilde{O} \left( \kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

iterations.

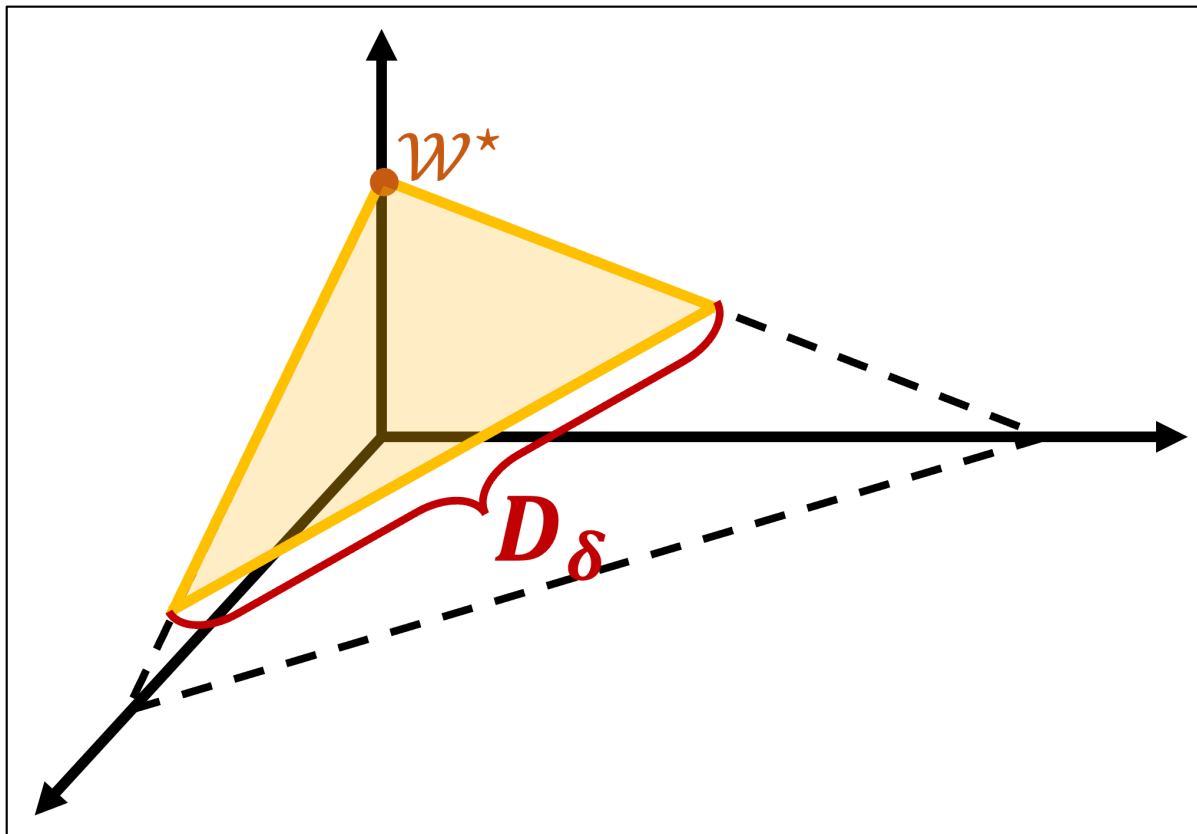
**Matrix** condition number of  $A$ :

$$\kappa := \sigma_{\max}^+(A) / \sigma_{\min}^+(A)$$

“Sublevel-set geometry”

$D_\delta$ : Diameter of  $\delta$ -sublevel set  $\mathcal{W}_\delta$

$$D_\delta := \max_{\bar{w}, \hat{w} \in \mathcal{W}_\delta} \|\bar{w} - \hat{w}\|$$

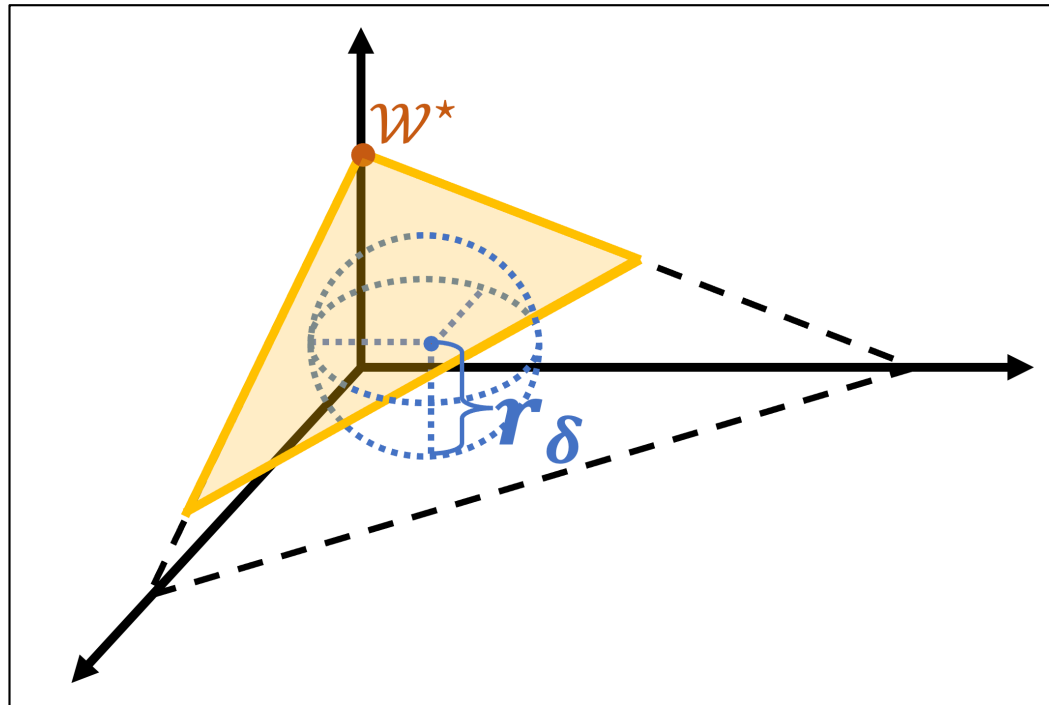


# $r_\delta$ : “Conic Radius” of $\mathcal{W}_\delta$

$$r_\delta := \max_{r \geq 0, w \in \mathcal{W}_\delta} r$$

s.t.  $B_w(r) \subset \mathcal{K} \times \mathcal{K}^*$

$r_\delta$  is the radius of the maximum ball inscribed in  $\mathcal{K} \times \mathcal{K}^*$  and centered at a point in  $\mathcal{W}_\delta$



# Target: $\varepsilon$ -optimal solution

$(x, s)$  is an  $\varepsilon$ -optimal solution if:

- distance to each type of constraint is no larger than  $\varepsilon$ , and
- the duality gap is not larger than  $\varepsilon$

$(x, s)$  is an  $\varepsilon$ -optimal solution if:

- $\text{Dist}(x, \{x | Ax = b\}) \leq \varepsilon$
- $\text{Dist}(x, \mathcal{K}) \leq \varepsilon$
- $\text{Dist}(s, \{s | \exists y \text{ s.t. } A^\top y + s = c\}) \leq \varepsilon$
- $\text{Dist}(s, \mathcal{K}^*) \leq \varepsilon$
- $c^\top x - b^\top (AA^\top)^{-1} A(c - s) \leq \varepsilon$

# Worst-case complexity of PDHG

**Theorem** [Xiong and F 2024]: Suppose  $w^*$  is unique. PDHG computes an  $\varepsilon$ -optimal solution within

$$\tilde{O} \left( \kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

iterations.

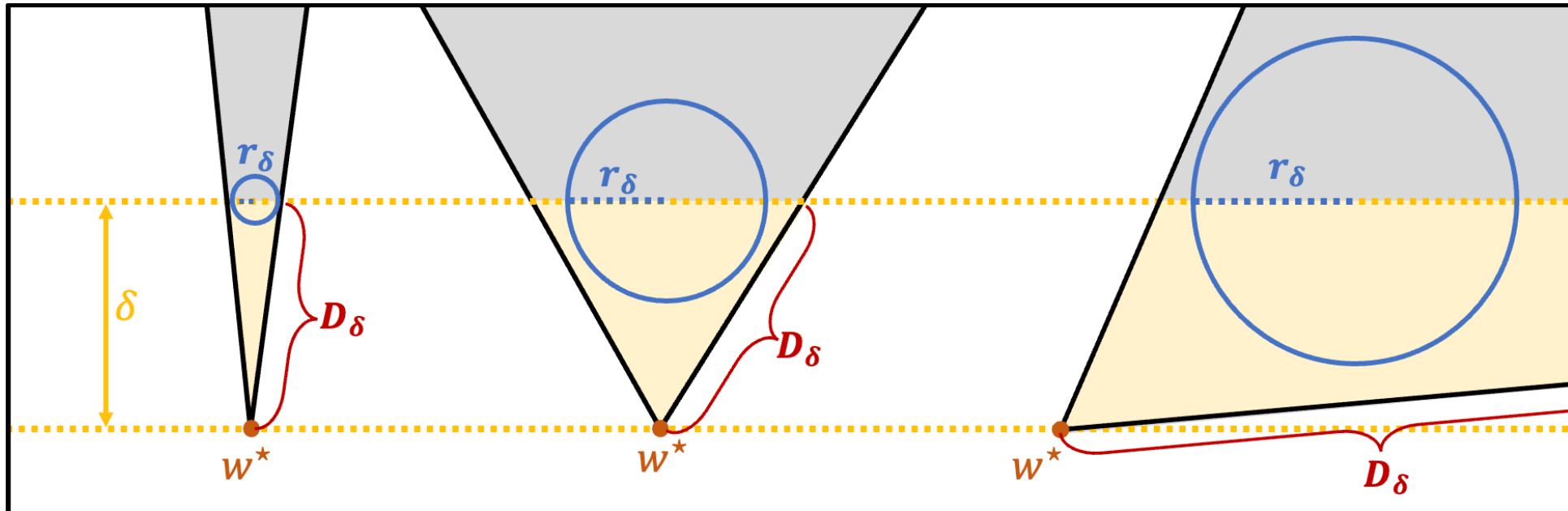
**Matrix** condition number of  $A$ :

$$\kappa := \sigma_{\max}^+(A) / \sigma_{\min}^+(A)$$

“Sublevel-set geometry”

# Local Geometry of and $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$ in the case of LP

When  $w^*$  is unique and  $\delta$  is sufficiently small,  $\mathcal{W}_\delta$  is a slice of a pointed cone at  $w^*$ .



Very small  $r_\delta$   
Intermediate  $D_\delta$



Intermediate  $r_\delta$   
Intermediate  $D_\delta$



Intermediate  $r_\delta$   
Very large  $D_\delta$

# Worst-case complexity of PDHG (under unique optima)

Matrix condition number of  $A$ :  $\kappa = \sigma_{\max}^+(A)/\sigma_{\min}^+(A)$

**Theorem** [Xiong and F 2024]: Suppose  $w^*$  is unique. PDHG computes an  $\varepsilon$ -optimal solution within

$$\tilde{O} \left( \kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

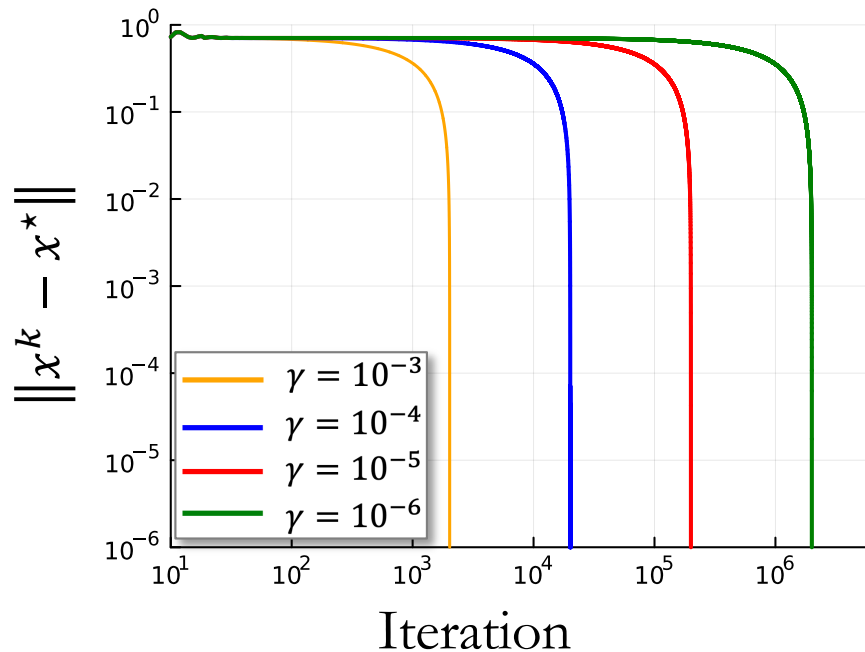
iterations.

Matrix condition number

Local geometric condition

# Recall the $P(\gamma)$ example

$$\begin{aligned} \min_{x_1, x_2} \quad & -(1 + \gamma)x_1 - x_2 \\ P(\gamma): \quad & \text{s.t. } x_1 + x_2 = 1 \\ & x_1 \geq 0, x_2 \geq 0 \end{aligned}$$

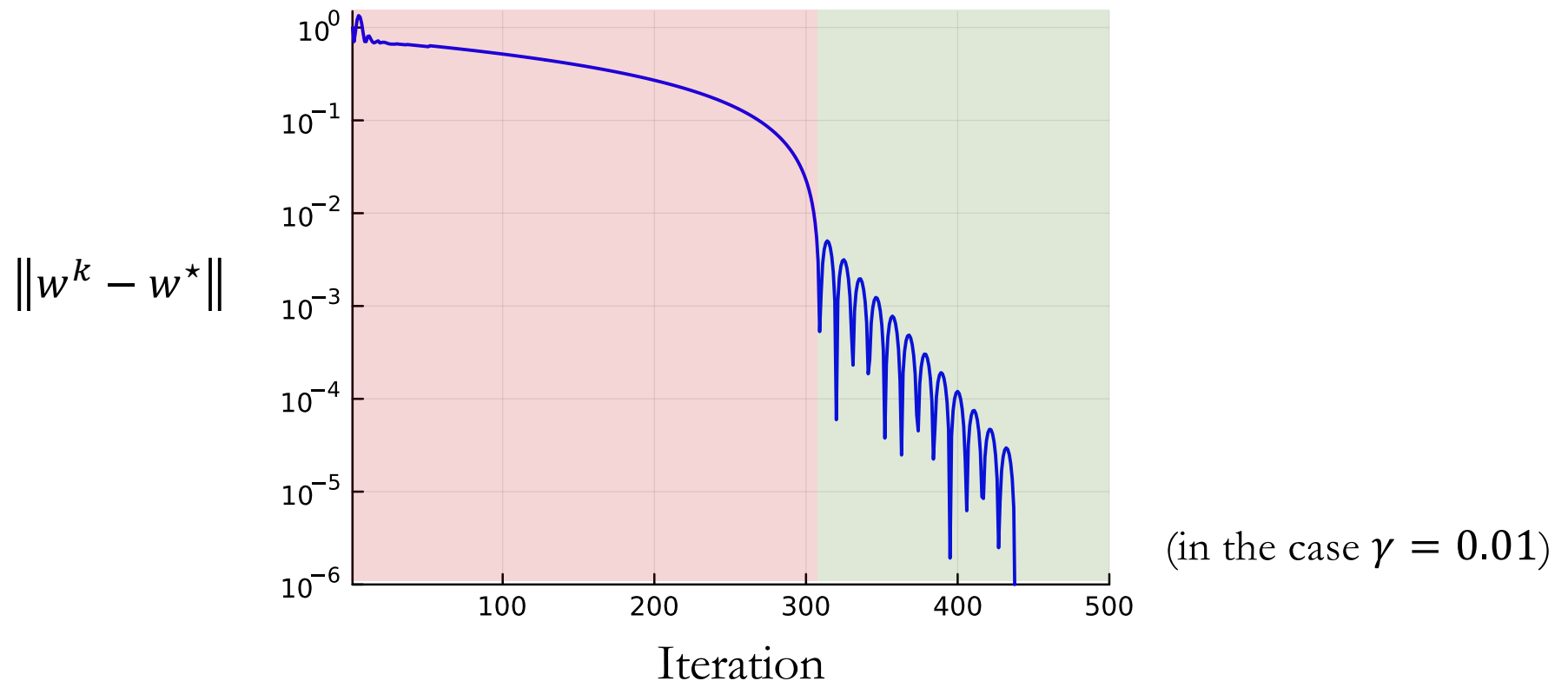


$\gamma$	Iteration count	$\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$
$10^{-3}$	2017	1414.2
$10^{-4}$	20023	14142.1
$10^{-5}$	200029	141421.4
$10^{-6}$	2000035	1414213.6

Observation: the number of iterations is linear in  $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$

# The Observed Two-Stage Performance of PDHG

## Typical convergence performance of PDHG



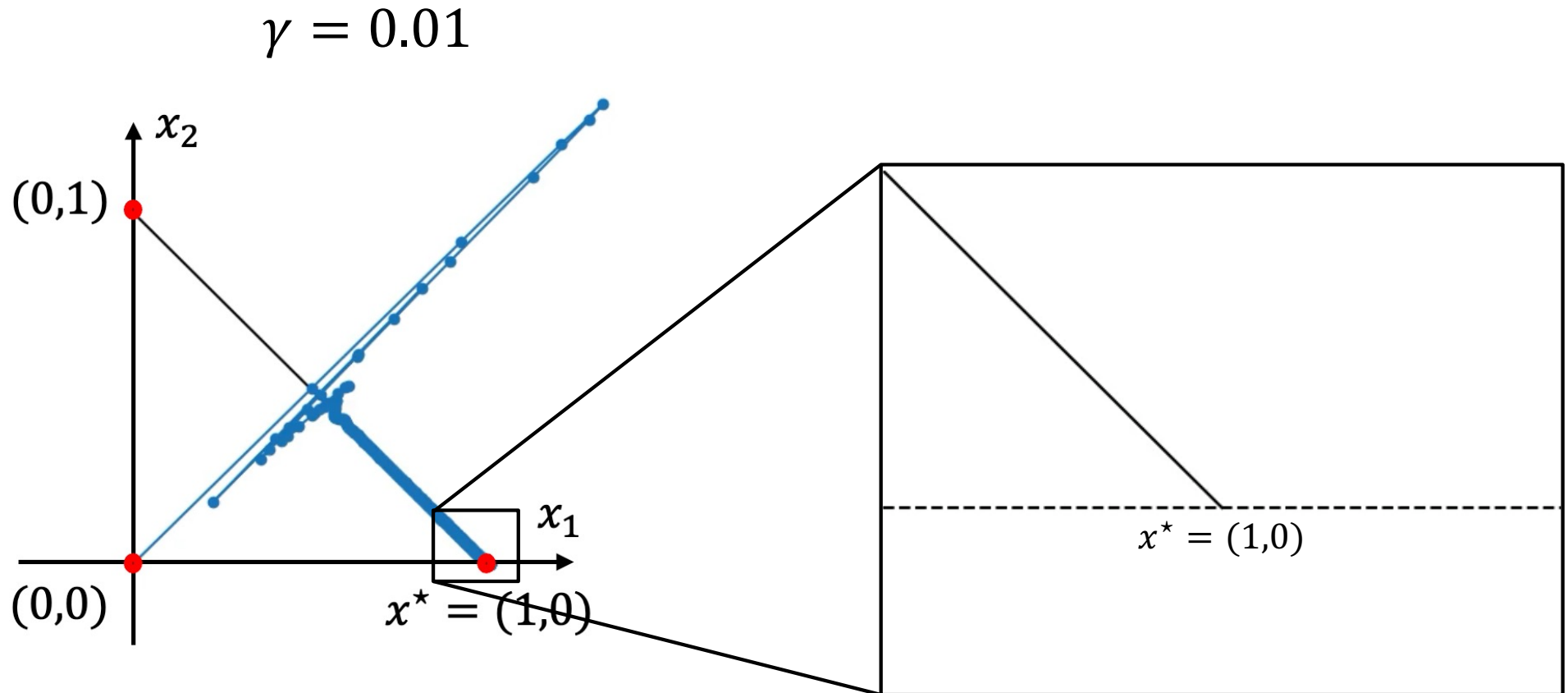
**Stage I:**

Finite-time basis identification

**Stage II:**

Fast local convergence

# The Observed Two-Stage Performance of PDHG



**Stage I:**

Finite-time basis identification

**Stage II:**

Fast local convergence

# Two-Stage Performance of PDHG

## Theorem (Stage I: Finite-time basis identification) [X, 2024a]:

The solution  $x^k$  identifies the optimal basis  $\mathcal{B}$  (i.e.  $\text{supp}(x^k) = \mathcal{B}$ ) for all  $k \geq T_{\text{basis}}$ , where:

$$T_{\text{basis}} := \tilde{O} \left( \kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \right)$$

## Theorem (Stage II: Fast local convergence) [X, 2024a]:

After identifying the optimal basis, PDHG computes an  $\varepsilon$ -optimal solution within additional

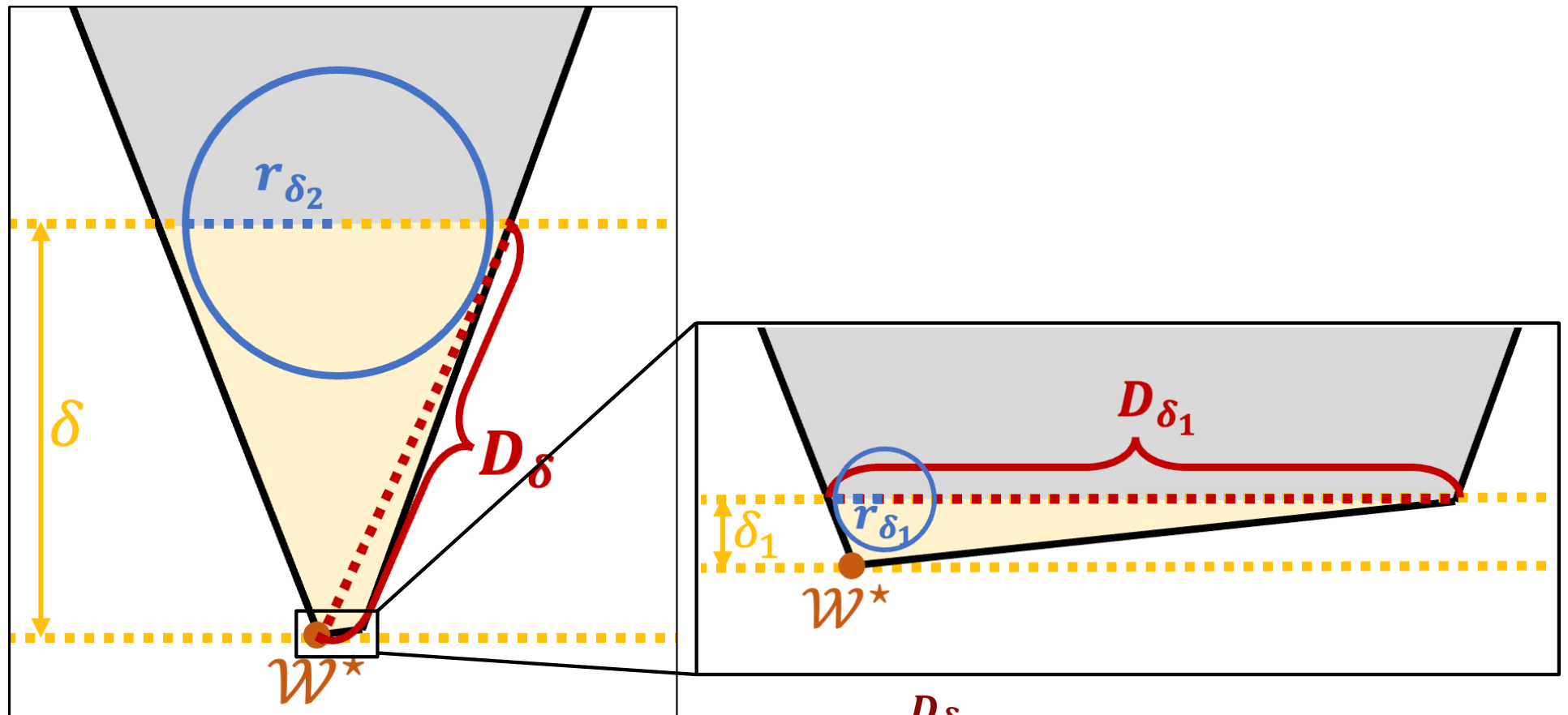
$$T_{\text{local}} := \tilde{O} \left( \|B^{-1}\| \|A\| \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

iterations.

$\|B^{-1}\| \|A\|$ : Only matrix condition numbers.

Stage II converges faster because it is essentially just solving a system of equations...

# Another example

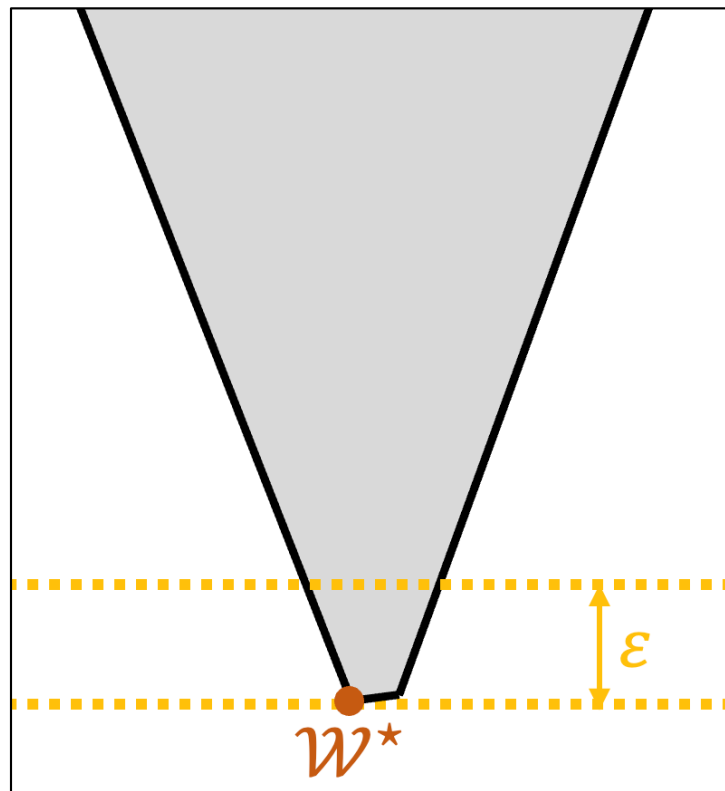


For  $\delta_2 > \delta_1$ ,  $\frac{D_{\delta_2}}{r_{\delta_2}}$  becomes smaller/better

$\frac{D_{\delta_1}}{r_{\delta_1}}$  is very large/bad  
(due to the small  $r_{\delta_1}$ )

Is  $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  the only geometric condition?

Suppose we want an  $\varepsilon$ -optimal solution:



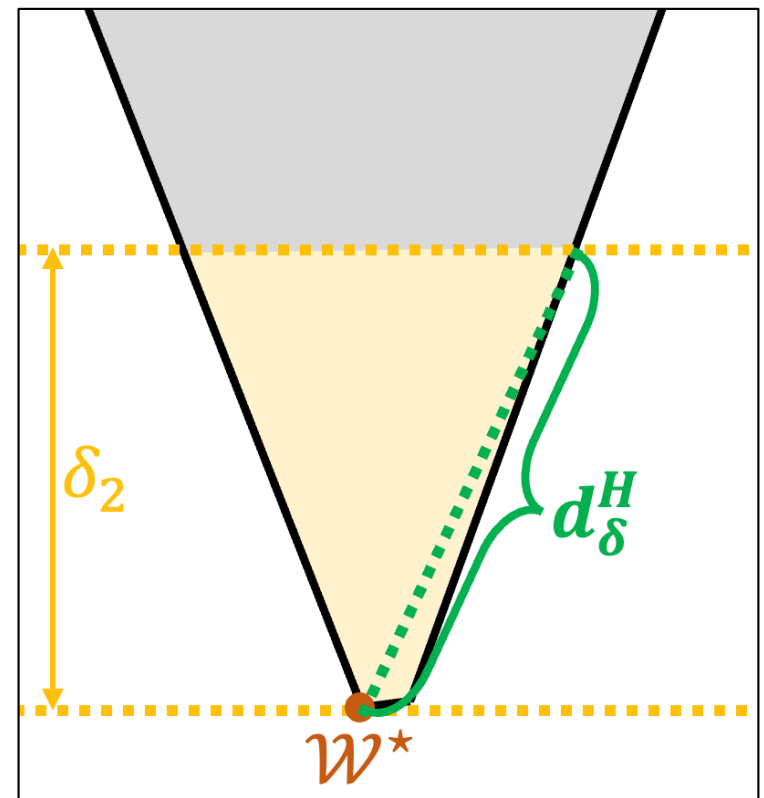
Intuition: The very-local bad geometry should not have a significant impact when the iterates of the algorithm have not yet reached the local neighborhood.

Is  $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  the only geometric condition?

We will need a third geometric measure that is designed to capture the notion of the level set “being close to  $\mathcal{W}^*$ ”

$$d_\delta^H := \max_{w \in \mathcal{W}_\delta} \text{Dist}(w, \mathcal{W}^*)$$

Hausdorff distance from  $\mathcal{W}_\delta$  to  $\mathcal{W}^*$



# Our General Conic Optimization Computational Guarantee

**Theorem** [Xiong and F 2024]: The number of PDHG iterations required to compute an  $\varepsilon$ -optimal solution is upper bounded by:

$$\tilde{O} \left( \kappa \cdot \max \left\{ \frac{D_\delta}{r_\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right), \frac{d_\delta^H}{\varepsilon} (1 + \text{Dist}(0, \mathcal{W}^*)) \right\} \right)$$

for each  $\delta > 0$ .

How good the geometry of  $\mathcal{W}_\delta$  is

How close  $\mathcal{W}_\delta$  is to  $\mathcal{W}^*$

**Implication:** If there is a  $\delta$ -sublevel set that

(i) has good geometry and (ii) is close to the optimal solution set, then PDHG may converge faster.

Remark: This result holds for LP with multiple optima, and for general conic optimization too.

$D_\delta$ : Diameter of ...

$r_\delta$ : Conic radius of ...

$d_\delta^H$ : Hausdorff distance ...

# Our General Conic Optimization Computational Guarantee

**Theorem** [Xiong and F 2024]: The number of PDHG iterations required to compute an  $\varepsilon$ -optimal solution is upper bounded by:

$$\tilde{O} \left( \kappa \cdot \max \left\{ \frac{D_\delta}{r_\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right), \frac{d_\delta^H}{\varepsilon} (1 + \text{Dist}(0, \mathcal{W}^*)) \right\} \right)$$

for each  $\delta > 0$ .

Small when  $\mathcal{W}_\delta$  has good geometry

Small when  $\mathcal{W}_\delta$  is close to  $\mathcal{W}^*$

Q: Can we improve  $\frac{D_\delta}{r_\delta}$  and  $d_\delta^H$  ?

**A: Yes, by using Hessian Rescaling**

$D_\delta$ : Diameter of ...

$r_\delta$ : Conic radius of ...

$d_\delta^H$ : Hausdorff distance ...

Consider the generic case of LP: LPs with unique optima

Matrix condition number of  $A$ :  $\kappa = \sigma_{\max}^+(A)/\sigma_{\min}^+(A)$

**Theorem** [Xiong and F 2024]: Suppose  $w^*$  is unique. PDHG computes an  $\varepsilon$ -optimal solution within

$$\tilde{O} \left( \kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

iterations.

Matrix condition number

Local geometric condition

[Xiong, 2024]:

- This bound has a closed-form expression
- PDHG has local fast convergence in a small neighborhood of  $w^*$
- This bound is  $\tilde{O} \left( n^{2.5} \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$  with high probability



# Establishing “Average-case” performance guarantees

# Probabilistic Analysis to shrink the gap between worst-case analysis and computational practice

	Worst-case complexity	Average-case complexity
Simplex Method	Exponentially poor [Klee and Minty, 1972]	$O(n^2)$ Borgwardt, Smale, Blair, Adler, Haimovich, Todd, Megiddo, et al. (~1980s)
Interior Point Method	$O(\sqrt{n}L)$ “Bit-length” $L$ might be huge (~1980s)	$O(\sqrt{n} \cdot \ln n)$ Todd, Ye, Anstreicher, Potra et al. (~1990s)
PDHG	Exponentially poor	$\tilde{O}(n^{2.5})$ with high probability [X, 2024b]

Note: I have not seen any other results showing that first-order methods for LP can achieve polynomial-time complexity with high probability.

# Iterations Bounds in Closed Form of the Optimal Basis/Solution

A formula for  $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  using the optimal basis/solution

### Strict Complementary Slackness (under unique optima):

Let  $\mathcal{B} := \text{supp}(x^*)$  and  $\mathcal{N} := \text{supp}(s^*)$ . Then  $(\mathcal{B}, \mathcal{N})$  is a partition of  $\{1, 2, \dots, n\}$ .

We use “ $\approx$ ” to denote being equivalent up to an absolute constant (2).

#### Lemma [Xiong 2024]

Under unique optima, let  $B := A_{\mathcal{B}}$ ,  $N := A_{\mathcal{N}}$ . Then

$$\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \approx (\|x^*\|_1 + \|s^*\|_1) \cdot \max \left\{ \max_{i \in [m]} \frac{\|(B^{-1}N)_{i,:}\| + 1}{x_{\mathcal{B}(i)}^*}, \max_{j \in [n-m]} \frac{\|(B^{-1}N)_{:,j}\| + 1}{s_{\mathcal{N}(j)}^*} \right\}.$$

$\ell_1$ -norm of the optimal solution

$$= \max_i \left( \frac{1 + \text{norm of the } i\text{-th row of } (B^{-1}N)}{i\text{-th component of } x_{\mathcal{B}}^*} \right)$$

$$= \max_j \left( \frac{1 + \text{norm of the } j\text{-th column of } (B^{-1}N)}{j\text{-th component of } s_{\mathcal{N}}^*} \right)$$

$B^{-1}A = (I, B^{-1}N)$  is the simplex method tableau

A formula for  $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  using the optimal basis/solution

### Strict Complementary Slackness:

Let  $\mathcal{B} := \text{supp}(x^*)$  and  $\mathcal{N} := \text{supp}(s^*)$ . Then  $(\mathcal{B}, \mathcal{N})$  is a partition of  $\{1, 2, \dots, n\}$ .

We use “ $\approx$ ” to denote being equivalent up to an absolute constant (2).

#### Lemma [Xiong 2024]

Under unique optima, let  $B := A_{\mathcal{B}}$ ,  $N := A_{\mathcal{N}}$ . Then

$$\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \approx (\|x^*\|_1 + \|s^*\|_1) \cdot \max \left\{ \max_{i \in [m]} \frac{\|(B^{-1}N)_{i,:}\| + 1}{x_{\mathcal{B}(i)}^*}, \max_{j \in [n-m]} \frac{\|(B^{-1}N)_{:,j}\| + 1}{s_{\mathcal{N}(j)}^*} \right\}.$$

Furthermore,  $\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  has a simple upper bound:

$$\lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \leq 2 \cdot \frac{\|x^*\|_1 + \|s^*\|_1}{\min_{i \in [n]} x_i^* + s_i^*} \cdot \|B^{-1}A\|$$

Ratio of  $\ell_1$ -norm to the smallest nonzero

Norm of the simplex tableau at the optimal solution

# An iteration bound using the optimal basis/solution

**Theorem** [Xiong 2024]: PDHG computes an  $\varepsilon$ -optimal solution within:

$$\tilde{O} \left( \kappa \cdot \frac{\|x^*\|_1 + \|s^*\|_1}{\min_{i \in [n]} x_i^* + s_i^*} \cdot \|B^{-1}A\| \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

iterations.

The smallest nonzero  $\left( \min_{i \in [n]} x_i^* + s_i^* \right)$  plays a key role in other methods as well:

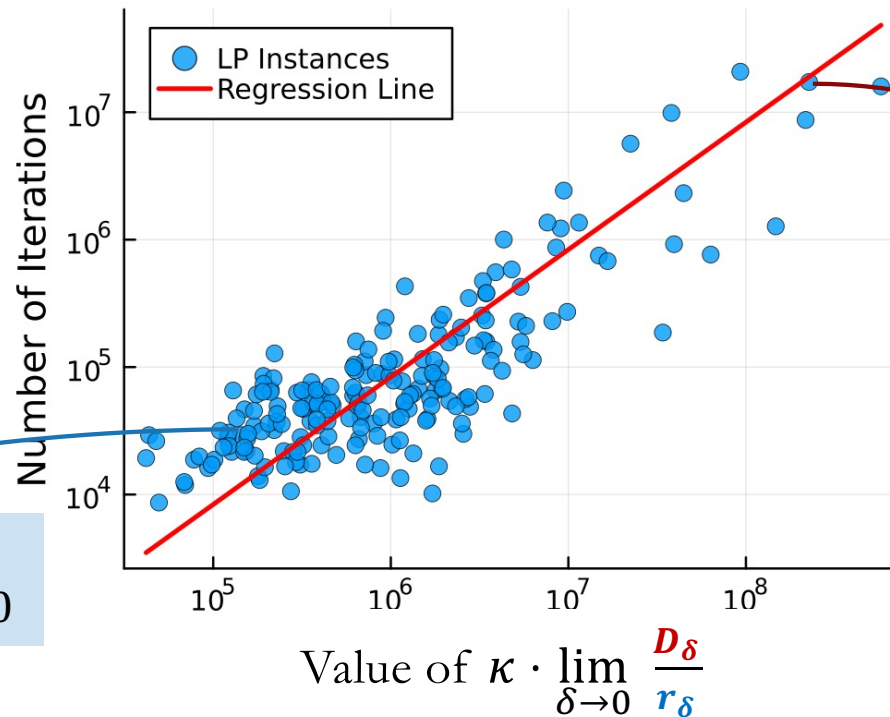
[Güler and Ye, 1993] [Ye, 1992] [Mehrotra and Ye, 1993] [Ye, 2011] ...

The product of  $\left( \frac{\|x^*\|_1 + \|s^*\|_1}{\min_{i \in [n]} x_i^* + s_i^*} \right)$  and a norm of  $B^{-1}A$  also appears in IPM complexity:

[Potra, 1994][Anstreicher, Ji, Potra and Ye, 1999] ...

# Validation Experiments on random LP instances

Iterations Required to solve LP instance to  $\varepsilon = 10^{-8}$



200 random LP instances  
with  $m = 60$ , and  $n = 120$

Regression line:  
 $\log_{10}(\text{PDHG Iteration})$   
 $= \log_{10} \kappa \Phi - 1.078$

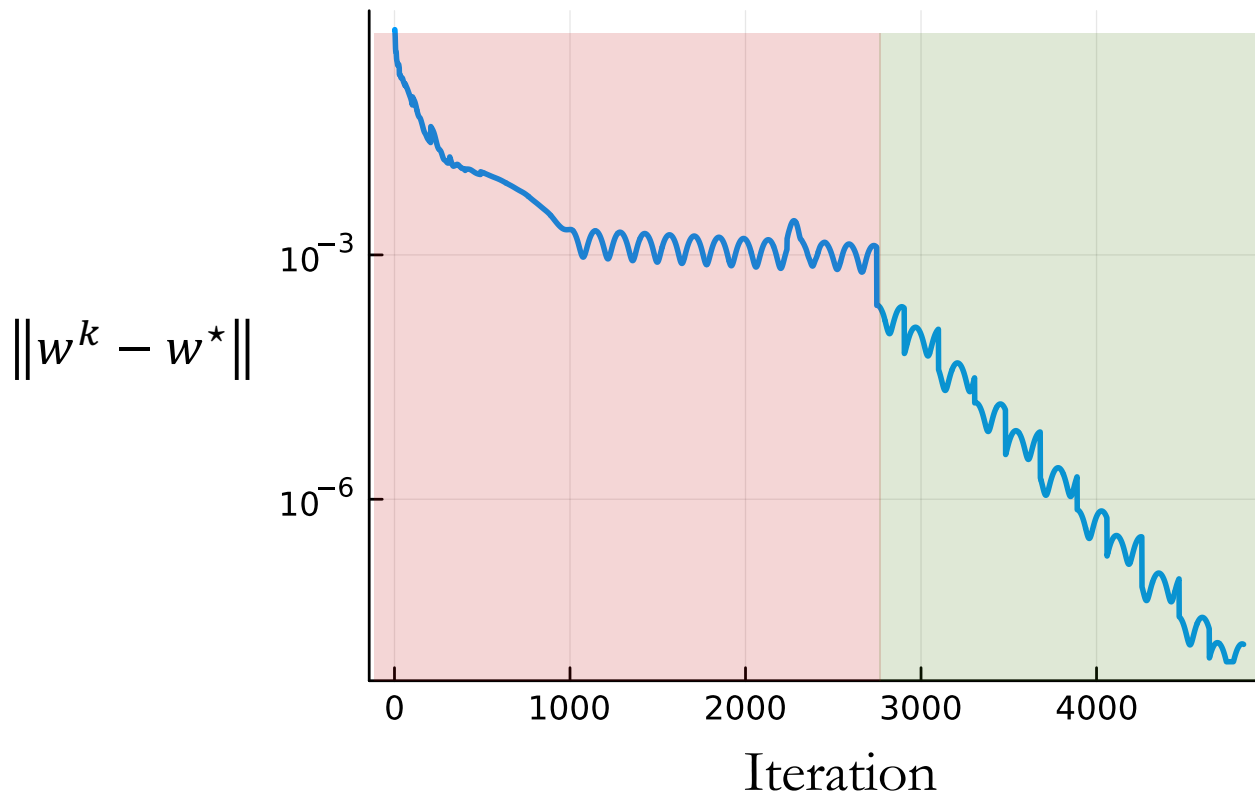
1. Fairly obvious linear dependence on  $\kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$
2. Although  $\kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  can be extremely large in practice (or not), the random LP instances typically have smaller values of  $\kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$ . More on this later...

Regarding the two-stage performance of PDHG on LP instances with unique optima

# Two-Stage Performance of PDHG

Two-stage performance of PDHG on general LP [Lu and Yang 2023]

Typical convergence performance of PDHG



**Stage I:**

Finite-time basis identification

**Stage II:**

Fast local convergence

# Two-Stage Performance of PDHG (under unique optima)

## **Theorem (Stage I: Finite-time basis identification)** [Xiong 2024]:

Let  $w^k = (x^k, s^k)$  denote the  $k$ -th iteration solution. The solution  $x^k$  identifies the optimal basis  $\mathcal{B}$  (i.e.  $\text{supp}(x^k) = \mathcal{B}$ ) for all  $k \geq T_{\text{basis}}$ , where :

$$T_{\text{basis}} := \tilde{O} \left( \kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta} \right)$$

## **Theorem (Stage II: Fast local convergence)** [Xiong 2024]:

Once the optimal basis has been determined, PDHG computes an  $\varepsilon$ -optimal solution within an additional

$$T_{\text{local}} := \tilde{O} \left( \|B^{-1}\| \|A\| \cdot \ln \left( \frac{1}{\varepsilon} \right) \right)$$

iterations.

$\|B^{-1}\| \|A\|$ : Only matrix condition numbers.

Stage II converges faster because it is not affected by the sublevel set geometry

Can we explain the practical performance of PDHG using high-probability complexity analysis?

# Todd's Classic Random LP Model

## Definition (Random LP Model):

Select  $\mathcal{B} \subset [n]$ , and  $|\mathcal{B}| = m$  and the solution  $(\hat{x}, \hat{s})$  is distributed as follows:

$$\begin{aligned}\hat{x}_{\mathcal{B}} &\sim |\mathcal{N}(0,1)|^m, & \hat{x}_{\mathcal{N}} &= 0, \\ \hat{s}_{\mathcal{B}} &= 0, & \hat{s}_{\mathcal{N}} &\sim |\mathcal{N}(0,1)|^{n-m},\end{aligned}$$

The random LP of the above optimal solution is distributed as follows:

$$A \sim \mathcal{N}(0,1)^{m \times n}, \quad b = A\hat{x}, \quad c = \hat{s}.$$

From a classic random LP model of [Todd, 1991].

Variants studied for IPM [Ye, 1994], [Anstreicher et al., 1999]

We use unit variance for simplicity of result.

This LP model has unique optima with probability = 1

# Two-Stage Performance of PDHG

## Theorem (Stage I: Finite-time basis identification) [Xiong 2024]:

Let  $T_{\text{basis}}$  denote the number of PDHG iterations to identify the optimal basis. Then it holds for any  $\delta \in \left(\frac{1}{2c_0n}, 1\right)$  that

$$\mathbb{P} \left[ T_{\text{basis}} \leq \tilde{O} \left( \frac{n^{2.5}}{\delta} \right) \right] \geq 1 - \delta.$$

## Theorem (Stage II: Fast local convergence) [Xiong 2024]:

After  $T_{\text{basis}}$  iterations, let  $T_{\text{local}}$  denote the number of additional PDHG iterations to compute an  $\epsilon$ -optimal solution. Then it holds for any  $\delta \in (0,1)$  that

$$\mathbb{P} \left[ T_{\text{local}} \leq \tilde{O} \left( \frac{n}{\delta} \cdot \ln \left( \frac{1}{\epsilon} \right) \right) \right] \geq 1 - \delta.$$

Faster local linear convergence indicated by the probabilistic analysis.

**Theorem (Stage I: Finite-time basis identification)** [Xiong 2024]:

Let  $T_{\text{basis}}$  denote the number of PDHG iterations to identify the optimal basis. Then it holds for any  $\delta \in \left(\frac{1}{2^{c_0 n}}, 1\right)$  that

$$\mathbb{P} \left[ T_{\text{basis}} \leq \tilde{O} \left( \frac{n^{2.5}}{\delta} \right) \right] \geq 1 - \delta.$$

**Theorem (Stage II: Fast local convergence)** [Xiong 2024]:

After  $T_{\text{basis}}$  iterations, let  $T_{\text{local}}$  denote the number of additional PDHG iterations to compute an  $\varepsilon$ -optimal solution. Then it holds for any  $\delta \in (0, 1)$  that

$$\mathbb{P} \left[ T_{\text{local}} \leq \tilde{O} \left( \frac{n}{\delta} \cdot \ln \left( \frac{1}{\varepsilon} \right) \right) \right] \geq 1 - \delta.$$

This provides a possible explanation for why PDHG works well in practice (polynomial-time in most cases), even though  $\kappa \cdot \lim_{\delta \rightarrow 0} \frac{D_\delta}{r_\delta}$  may take extreme values in the worst case.

But PDHG bound has a heavier tail compared with the two classic methods:

- PDHG has polynomial high-probability complexity
- IPM has polynomial average-case complexity. [Anstreicher, et al., 1999]

# Summary/Remarks

- The convergence rate of PDHG for conic optimization is related to the geometry of primal-dual sublevel sets measured with  $D_\delta, r_\delta, d_\delta^H$
- For LP instances with unique optima, the iteration bound has a closed-form expression
- For LP instances with unique optima, PDHG has faster local convergence after identifying the optimal basis
- PDHG is polynomial-time with high probability

## Remark:

- These results relied only on PDHG's average iterate convergence and non-expansiveness properties. Similar results might also hold for other FOMs, in particular ADMM, EGM, ...



Thank you!